# Agency Theory: A Reading

## Prof Dr. Adler Haymans Manurung, CIFM, CIGS
## Dr. M. Jhonni Sinaga, CIPFM

# Agency Theory:
## A Reading

**Prof. Dr. Adler Haymans Manurung**
CIFM, CIMA, CIERM, CIABV, CMA, CBV, CERA, CIBG, CIGS, CIQnR, CIQaR

**Dr. M. Jhonni Sinaga, SE., MM.,** CIPFM

**Universitas Bhayangkara Jakarta Raya**

# Agency Theory:
## A Reading

@**Prof. Dr. Adler Haymans Manurung**
CIFM, CIMA, CIERM, CIABV, CMA, CBV, CERA, CIBG, CIGS, CIQnR, CIQaR
@**Dr. M. Jhonni Sinaga, S.E., M.M.,** CIPFM, CIERM

**Printed by :**
**The contents are out of the printing house's responsibility.**

I dedicate this book to my beloved wife, daughter, and son ; Marsaurina Yudiciana Sitanggang, Castelia Romauli, and Adry Gracio.

**Prof. Dr. Adler Haymans Manurung, SE., SH., M.Com., ME.**
CIFM, CMA, CIMA, CIERM, CIBG, CIABV, CIQnR,. CIQaR, CIRR

This book is dedicated to my beloved wife and both of my children; Rosalia Manurung, Calvin Jhon Junior, and Jessie Jhon Junior.

**Dr. M. Jhonni Sinaga,S.E., M.M.,** CIPFM, CIERM

# Foreword

Agency theory has gained its prominence in international organizations, academics, professional practices, and corporate bodies over the years now. An agency relationship arises when a provider of funds appoints another to manage his interest. Proponents of agency theory believe that there is a tendency for agents, when left unmonitored, to engage in self-interest activities to the principal's detriment. Meanwhile, the degree of board independence is positively and significantly related to firm performance, especially in government-controlled firms and in firms with lower information acquisition and monitoring costs. The three phenomena require the presence of management control system to meet the interests of fund providers.

Agency theory posits different organizational, behavioral, economical and controlling roles, and it is a potent framework which can be extricated in promulgation of the management control systems. Furthermore, the implementation of a control mechanism depends on the amount and contents of the public and/ or private information that exist in the domain of the managerial accounting system. Furthermore, the implementation of a control mechanism depends on the amount and contents of the public and/ or private information that exist in the domain of the managerial accounting system.

Indeed, what makes CEO pay both interesting and complicated is the fact that the efficient contracting, managerial power, and political paradigms co-exist and interact. In introducing plans that tie pay more strongly to performance as demanded by shareholders, directors routinely agree to pay more than necessary to compensate for the increased risk.

This book is prepared for everyone that has strong intention to empower his or her theories and empirical information dealing with agency theory which is much helpful to make decisions in regards meeting the interests of fund providers. It is also a reading material for students attending doctoral program in the field of finance. Various international journals, both theoretical and empirical, were collected in the context of teaching materials.

Constructive criticisms are really needed to make this book more useful and in good quality. This book is going to be continuously and accordingly updated as well.

Sincerely,

Prof. Dr. Adler Haymans Manurung, CIFM, CMA

Dr. Jhonni Sinaga CIPFM

# LIST OF CONTENT

# The Economic Theory of Agency: The Principal's Problem

*By* STEPHEN A. ROSS*

The relationship of agency is one of the oldest and commonest codified modes of social interaction. We will say that an agency relationship has arisen between two (or more) parties when one, designated as the agent, acts for, on behalf of, or as representative for the other, designated the principal, in a particular domain of decision problems. Examples of agency are universal. Essentially all contractural arrangements, as between employer and employee or the state and the governed, for example, contain important elements of agency. In addition, without explicitly studying the agency relationship, much of the economic literature on problems of moral hazard (see K. J. Arrow) is concerned with problems raised by agency. In a general equilibrium context the study of information flows (see J. Marschak and R. Radner) or of financial intermediaries in monetary models is also an example of agency theory.

The canonical agency problem can be posed as follows. Assume that both the agent and the principal possess state independent von Neumann-Morgenstern utility functions, $G(\cdot)$ and $U(\cdot)$ respectively, and that they act so as to maximize their expected utility. The problems of agency are really most interesting when seen as involving choice under uncertainty and this is the view we will adopt. The agent may choose an act, $a \in A$, a feasible action space, and the random payoff from

this act, $w(a, \theta)$, will depend on the random state of nature $\theta(\epsilon\Omega$ the state space set), unknown to the agent when $a$ is chosen. By assumption the agent and the principal have agreed upon a fee schedule $f$ to be paid to the agent for his services. The fee, $f$, is generally a function of both the state of the world, $\theta$, and the action, $a$, but we will assume that the action can influence the parties and, hence, the fee only through its impact on the payoff. This permits us to write,

$$(1) \qquad f = f(w(a, \theta); \theta).$$

Two points deserve mention. Obviously the choice of a fee schedule is the outcome of a bargaining problem or, in large games, of a market process. Much of what we have to say is relevant for this view but we will not treat the bargaining problem explicitly. Second, while it is possible to conceive of the fee as being directly functionally dependent on the act, the theory loses much of its interest, since without further conditions, such a fee can always be chosen as a Dirac $\delta$-function forcing a particular act (see S. Ross). In some sense, then, we are assuming that only the payoff is operational and we will take this point up below. Now, the agent will choose an act, $a$, so as to

$$(2) \qquad \max_{a} E_{\theta}\{G[f(w(a, \theta); \theta)]\},$$

where the agent takes the expectation over his subjectively held probability distribution. The solution to the agent's problem involves the choice of an optimal act, $a^o$, conditional on the particular fee schedule, i.e., $a^o = a(\langle f \rangle)$, where $a(\cdot)$ is a

1

mapping from the space of fee schedules into $A$.

If the principal has complete information about the fee to act mapping, $a(\langle f \rangle)$, he will now choose a fee so as to

$$(3) \quad \max_{\langle f \rangle} E_{\theta}\{ U[w(a(\langle f \rangle)), \theta) \\ \qquad - f(w(a(\langle f \rangle), \theta); \theta)] \},$$

where the expectation is taken over the principal's subjective probability distribution over states of nature. If the principal is not fully informed about $a(\cdot)$, then $a(\cdot)$ will be a random function from his point of view. Formally, at least, by appropriately augmenting the state space the criterion (3) could still be made to apply. In general some side constraints on $\langle f \rangle$ would also have to be imposed to insure that the problem possesses a solution (see Ross). A market-imposed minimum expected fee or expected utility of fee by the agent would be one economically sensible constraint:

$$(4) \quad E_{\theta}\{ G[f(w(a, \theta); \theta)] \} \geq k.$$

Since utility functions are assumed to be independent of states, $\theta$, one of the important reasons for a fee to depend directly on $\theta$ would be if individual subjective probability distributions differed. In what follows we will assume that both the agent and the principal share the same subjective beliefs about the occurrence of $\theta$ and write the fee as a function of the payoff only,

$$(5) \quad f = f(w(a, \theta)).$$

Notice that this interpretation would not in general be permissible if the principal lacked perfect knowledge of $a(\cdot)$. More importantly, though, surely aside from simple comparative advantage, for some questions the *raison d'être* for an agency relationship is that the agent (or the principal) may possess different (better or finer) information about the states of

the world than the principal (agent). If we abstract from this possibility we will have to show that we are not throwing out the baby with the bath water.

Under this assumption the problem is considerably simplified but much of interest does remain. Suppose, first, that we are simply interested in the properties of Pareto-efficient arrangements that the agent and the principal will strike. Notice that the optimal fee schedule as seen by the principal is found by solving (3) and is dependent on the desire to motivate the agent. In general, then, we would expect such an arrangement to be Pareto-inefficient, but we will return to this point below. The family of Pareto-efficient fee schedules can be characterized by assuming that the principal and the agent cooperate to choose a schedule that maximizes a weighted sum of utilities

$$(6) \quad \max_{\langle f \rangle} E\{ U[w - f] + \lambda G[f] \},$$

where $\lambda$ is a relative weighting factor (and where strategies have been randomized to insure convexity). K. Borch recognized that the solution to (6) is obtained by maximizing the function internal to the expectation which requires setting

$$(\text{P.E.}) \quad U'[w - f] = \lambda G'[f]$$

when $U$ and $G$ are monotone and concave. (See H. Raiffa for a good exposition.) The P.E. condition defines the fee schedule, $f(\cdot)$, as a function of the payoff $w$ (and the weight, $\lambda$). (See R. Wilson (1968) or Ross for a fuller discussion of this derivation and the functional aspect of the fee schedule.)

An alternative approach to finding optimal fee schedules was first proposed by Wilson in the theory of syndicates and studied by Wilson (1968, 1969) and Ross. This is the similarity condition that solves for the fee schedule by setting

(S)          $U[w - f] = aG[f] + b$

for constants $a > 0$, $b$. If $\langle f \rangle$ satisfies $S$ then, given the fee schedule, it should be clear that the agent and the principal have identical attitudes towards risky payoffs and, consequently, the agent will always choose the act that the principal most desires. Ross was able to completely characterize the class of utility functions that satisfied both P.E. and $S$ (for a range of $\lambda$) and show that in such situations the fee schedule is (affine) linear, $L$, in the payoff. (The class is simply that of pairs $\langle U, G \rangle$ with linear risk tolerance,

$$-\frac{U'}{U''} = cw + d \quad \text{and} \quad -\frac{G'}{G''} = cw + e,$$

where $c$, $d$ and $e$ are constants.) In fact, it can be shown that any two of $S$, P.E., or $L$ imply the third.

A question of interest that naturally arises is that of the relation that $S$ and P.E. bear to the exact solution to the principal's problem. (A comparable "agent's problem" can also be posed but we will not be concerned with that here. Some observations on such a problem are contained in Ross.) The solution to the principal's problem (3) subject to the constraint (4) and to the constraint imposed by the condition that the agent chooses the optimal act from his problem (2) can, under some circumstances, be posed as a classical variational problem. To do so we will assume that the payoff function is (twice) differentiable and that the agent chooses an optimal act, given a fee schedule, by the first order condition

(7)          $E_{\theta}\{G'[f(w)]f'(w)w_a\} = 0,$

where a subscript indicates partial differentiation. The principal's problem is now to

(8)          $\max_{\langle f \rangle} E_{\theta}\{H\} \equiv \max_{\langle f \rangle} E_{\theta}\{U[w - f]$

$$+ \Psi G'f'w_a + \lambda G\}$$

where $\Psi$ and $\lambda$ are Lagrange multipliers associated with the constraints (7) and (4) respectively. Changing variables to $V(\theta) \equiv f(w(a, \theta))$ where we have suppressed the impact of $a$ on $V$ and assuming, without loss of generality, that $\theta$ is uniformly distributed on $[0, 1]$ permits us to solve (8) by the Euler-Lagrange equation. Thus, at an optimum

(9)          $\dfrac{d}{d\theta}\left\{\dfrac{\partial H}{\partial V'}\right\} - \dfrac{\partial H}{\partial V}$

$$= U' + \Psi G'\frac{d}{d\theta}\left[\frac{w_a}{w_\theta}\right] - \lambda G' = 0;$$

or the marginal rate of substitution,

(10)          $\dfrac{U'}{G'} = \lambda - \Psi \dfrac{d}{d\theta}\left[\dfrac{w_a}{w_\theta}\right].$

This is an intuitively appealing result; the marginal rate of substitution is set equal to a constant as in the P.E. condition plus an additional term which captures the constraint (7) imposed on the principal by the need to motivate the agent. To determine the optimal act, $a$, we differentiate (8) with respect to $a$ which yields

(11)          $E_{\theta}\{U'[1 - f']w_a + \Psi G''(f'w_a)^2$

$$+ \Psi G'f''(w_a)^2 + \Psi G'f'w_{aa}\} = 0,$$

where we have made use of (7). Substituting the boundary conditions permits us to solve for the multipliers $\Psi$ and $\lambda$.

Like $S$ or P.E. (10) defines the fee schedule as a function of $w$. (Notice that we are tacitly assuming that, at least for the optimal act, the payoff is (*a.e.* locally) state invertible. This allows the fee to take the form of (5).) It follows that (10) will coincide with P.E. if and only if $\Psi$ is zero, or if $\Psi \neq 0$, we must have

3

$$(12) \qquad \frac{d}{d\theta}\begin{bmatrix} w_a \\ \frac{}{w_\theta} \end{bmatrix} = b(a),$$

a function of $a$ alone.

In particular, using these conditions we can ask what class of (pairs of) utility functions $\langle U, G \rangle$ has the property that, for any payoff structure, $w(a, \theta)$, the solution to the principal's problem is Pareto-efficient. Conversely, we can ask what class of payoff structures has the property that the principal's problem yields a Pareto-efficient solution for any pair of utility functions $\langle U, G \rangle$.

A little reflection reveals that the only pairs of $\langle U, G \rangle$ that could possibly belong to the first class must be those which satisfy S and P.E. for a range of schedules (indexed by the $\lambda$ weight in P.E.). Clearly if (10) is to be equivalent to P.E. for all payoff functions, $w (a, \theta)$, then $\Psi$ must be zero and the motivational constraint (7) must not be binding. For this to be the case, for an interval of values of $k$ (in (4)), the satisfaction of P.E. must imply that the agent chooses the principal's most desired act by (7). For any fee schedule, $\langle f \rangle$, the principal wants the act to be chosen to maximize $E_\theta\{U[w-f]\}$ which implies that

$$(13) \qquad E_\theta\{U'(1 - f')w_a\} = 0.$$

If (13) is to be equivalent to the motivational constraint (7) for all possible payoff structures, then we must have

$$(14) \qquad U'(1 - f') = G'f'$$

which, with P.E. (or (10) with $\Psi=0$) yields a linear fee schedule in the payoff. But, as shown in Ross, linearity of the fee schedule and P.E. imply the satisfaction of $S$ and the $\langle U, G \rangle$ pair must belong to the linear risk-tolerance class of utility functions described above.

Since the linear risk-tolerance class, while important, is very limited, we turn

now to the converse question of what payoff structures permit a Pareto-efficient solution for all $\langle U, G \rangle$ pairs. If $\Psi = 0$ we must, as before, have that the motivational constraint is not binding for all $\langle U, G \rangle$ or (13) must always imply (7). The implication will always hold if there exists an $a^*$ such that for all $a$ there is some choice of the state domain, $I$, for which

$$(15) \qquad w(a^*, \theta) \geq w(a, \theta), \qquad \theta \in I.$$

Conversely, from P.E., we must have that for all $G(\cdot)$

$$(16) \qquad E_\theta\{G'[f](1 - f')w_a\} = 0$$

implies (7) where $f$ is determined by P.E. Since $\langle U, G \rangle$ can always be chosen so as to attain any desired weightings of $w_a$ in (7) and (16) the special case of (15) is the only one for which motivation is irrelevant. Given (15) all individuals have a uniquely optimal act irrespective of their attitudes towards risk.

If $\Psi \neq 0$, then to assure Pareto efficiency we must satisfy (12). This is a partial differential equation and its solution is given by

$$(17) \qquad w(a, \theta) = H[\theta B(a) - C(a)],$$

where $H(\cdot)$, $B(\cdot)$ and $C(\cdot)$ are arbitrary functions. (The detailed computations are carried out in an appendix.) This is a rich and interesting class of payoff functions. In particular, (17) is a generalization of the class of functions of the form $l(\theta - a)$, where the object is to pick an act, $a$, so as to best guess the state $\theta$. It therefore includes, for example, traditional estimation problems, problems with a quadratic payoff function, and all problems with payoff functions of the form $|\theta - a|^\xi h(a)$, and many asymmetric ones as well. It is not, however, difficult to find plausible payoff functions which do not take the form of (17). (The class of the form (15) will generate such functions.)

We may conclude, then, that the class of payoff structures that simultaneously solve the principal's problem and lead to Pareto efficiency for all $\langle U, G \rangle$ pairs is quite important and quite likely to arise in practice.

In general, though, it is clear that the solution to the principal's problem will not be Pareto-efficient. This is, however, a somewhat naive view to take. Pareto efficiency as defined above assumes that perfect information is held by the participants. In fact, the optimal solution to the principal's problem implied that the fee-to-act mapping induced by the agent was completely known to the principal. In such a case it might be thought that the principal could simply tell the agent to perform a particular act. The difficulty arises in monitoring the act that the agent chooses. Michael Spence and Richard Zeckhauser have examined this problem in detail in the case of insurance. In addition, if agents are numerous the fee may be the only communication mechanism. While it might in principle be feasible to monitor the agent's actions, it would not be economically viable to do so.

The format of this paper has been such as to allow us to only touch on what is surely the most challenging aspect of agency theory; embedding it in a general equilibrium market context. Much is to be learned from such attempts. One would naturally expect a market to arise in the services of agents. Furthermore, in some sense, such a market serves as a surrogate for a market in the information possessed by agents. To the extent to which this occurs, the study of agency in market contexts should shed some light on the economics of information. To mention one more path of interest—in a world of true uncertainty where adequate contingent markets do not exist, the manager of the firm is essentially an agent of the shareholders. It can, therefore, be expected that

an understanding of the agency relationship will aid our understanding of this difficult question.

The results obtained here provide some of the micro foundations for such studies. We have shown that, for an interesting class of utility functions and for a very broad and relevant class of payoff structures, the need to motivate agents does not conflict with the attainment of Pareto efficiency. At the least, a callous observer might view these results as providing some solace to those engaged in econometric activity.

### APPENDIX

This appendix solves the partial differential equation (12) in the text. Integrating (12) over $\theta$ yields

$$\frac{\partial w}{\partial a} + \left[ b(a)\theta + c(a) \right] \frac{\partial w}{\partial \theta} = 0.$$

Along a locus of constant $w$,

$$\frac{d\theta}{da} = -\frac{\partial w/\partial a}{\partial w/\partial \theta} = b(a)\theta + c(a),$$

is a first order Bernoulli equation that integrates to

$$\theta = e^{\int b(a)} \left[ \int e^{-\int b(a)} c(a) + k \right],$$

where $k$ is a constant of integration. It follows that

$$w(a, \theta) = H[\theta B(a) - C(a)],$$

where

$$B(a) \equiv e^{-\int b(a)}$$

and

$$C(a) \equiv \int e^{-\int b(a)} c(a) + k,$$

and $H(\cdot)$ is an arbitrary function.

## REFERENCES

K. J. Arrow, *Essays in the Theory of Risk-Bearing,* Chicago 1970.

K. Borch, "Equilibrium in a Reinsurance Market," *Econometrica,* July 1962, *30,* 424–444.

J. Marschak and R. Radner, *The Economic Theory of Teams,* New Haven and London 1972.

H. Raiffa, *Decision Analysis; Introductory Lectures on Choices Under Uncertainty,* Reading, Mass. 1968.

S. Ross, "On the Economic Theory of Agency: The Principle of Similarity," *Proceedings of the NBER-NSF Conference on Decision Making and Uncertainty,* forthcoming.

M. Spence and R. Zeckhauser, "Insurance, Information and Individual Action," *Amer. Econ. Rev. Proc.,* May 1971, *61,* 380–387.

R. Wilson, "On the Theory of Syndicates," *Econometrica,* Jan. 1968, *36,* 119–132.

———, "The Structure of Incentives for Decentralization Under Uncertainty," *La Decision,* Editions Du Centre National De Le Recherche Scientifique, Paris 1969.

# Articles ◣3

# THE THEORY OF AGENCY:

## THE POLICING "PARADOX"

## AND REGULATORY BEHAVIOR

### Barry M. Mitnick ★

Relations of agency, or "acting for," are pervasive in complex societies. Examples include the worker-boss, physician-patient, adviser-administrator, and parent-child relations. Functional dependencies, among other reasons, determine that agency relationships will be extremely common. A key problem that principals in such relations face is that of insuring that the agent does in fact act for the principal. This paper will present a model of agency with policing and apply it to discuss regulatory behavior.

The policing "paradox" to be described is only an apparent, ascribed paradox, evident mostly in the short-run, and due generally to information effects. It occurs to observers of policing because of lack of knowledge about the goal state and information level of policing participants, and about the dynamic characteristics of the policing process. By identifying it we wish to emphasize policing as a process, with definite stages, rather than compress the argument in the way the usual holistic economic approach would proceed.

The policing arguments are developed in the context of the theory of agency, presently under development (see Mitnick 1973, 1974, 1975; for the economic theory of agency, see Ross 1973, 1974; see also Goldberg 1973, 1974; for a related model of policing in the context of the firm, see Alchian and Demsetz 1972), and managerial discretion models (see, e.g., Williamson 1964; Migue and Belanger 1974).

We have chosen a fiduciary function representation (see below), derived from Williamson (1964), rather than the budget-output line employed by Migue and

*Assistant Professor of Public Administration and Political Science, Ohio State University.

Belanger (1974) to depict discretionary resources. This is partly because we seek a general agency model of policing, wishing to examine returns to agent and principal directly rather than use the surrogate "output" for principal returns, and partly in order to be able to compare agent and principal returns directly. Both Williamson and Migue and Belanger have noted, as we do, that income and substitution effects may occur in managerial discretion models. The sources and circumstances of these effects in each case, however, are different. Williamson (1964) describes the effects of changes in tax rates in each of his models. Migue and Belanger note the effects of changes in demand, cost, and demand and cost elasticities on the budget-output function. The present paper discusses policing effects using the fiduciary function representation. (See also Mitnick 1973, 1974.)

We begin by describing fiduciary relations in the context of agency.

## I. Agency and Fiduciary Relations

The agent holding the *fiduciary norm* must act diligently, with the skills at his disposal, for the principal's goal, without regard for any other goals that may bear on his relation with the principal, including any self-goals.[1] The norm may be expected in contractual discretionary agency, typically under conditions of trust, under principal dependency, or under agent domination of the principal's interests. The *fiduciary function*, relation (in the mathematical sense), or set will be said to be the set of collections of specifications of the returns to each of the agent's goals, self- and other-, that are subject to the possible choice of the agent, given his contractual arrangement with the principal. We may also term this set the *institution function*, relation, or set since the possibly policed formal norms or rules of an institution,[2] perhaps an organization, may constrain the set of selections available to the agent. The *pure fiduciary* will choose that collection of specifications which contains the highest returns to his other-goals, regardless of the level of return to his self-goals. The *lexicographic*, or *"lexical" fiduciary* will after choosing the collections with the highest return to the principal, choose that collection which also has the highest return to his self-goals. These types may be contrasted with types of agents who choose for their self-goal first: The *pure*

---

[1]See Mitnick (1975). The arguments below regarding the policing "paradox" and regulatory behavior are from Mitnick (1974). Note that the usage of "fiduciary" is different from that employed by Curry and Wade (1968), who use it for entrepreneurs conceived as agents for groups, and do not require the restriction on consideration of competing, including self-, interests. General usage (see e.g., Riker 1962, pp. 24-28; Pitkin 1967; Seavey 1964) agrees with the sense here employed.

A *self-goal* is an objective of self-regarding preferences; an *other-goal* is an objective of other-regarding preferences. Here "self-regarding preferences" relate to concerns that are private, personal, egoistic, selfish, self-bettering, and so on; "other-regarding preferences" relate to such concerns of the "other" party.

[2]Blake and Davis (1964, p. 464) note that "norms clustered around a given functional requirement are often collectively designated as 'institution'."

*self-interest agent* chooses that collection of specifications which contains the highest returns to his self-goals, regardless of the level of return to his other-goals. The *lexical self-interest agent* chooses, first, those collections of specifications which contain the highest returns to his self-goals and, second, from that group of collections, that collection which has the highest return to his other-goals.[3] We posit that policing may involve a general developmental pattern (allowing skips in stages) from the pure self-interest agent to the lexical self-interest agent to the lexical fiduciary and finally to the pure fiduciary.

## II. Diagrammatic Exposition of Agency with Policing

Assume that the agent and principal are rational, have preferences that reflect "greed," i.e., "more is always better," and that the agent possesses a genuinely preferred self-goal and a genuinely preferred other-goal (the goal of the principal). Assume that the agent's indifference curves between "a" (return to self-goal as measured in "resources" devoted to it) and "p" (return to other-goal as measured in "resources" devoted to it) may be constructed, are continuous and that they have negative slope. Assume also that the slopes of the indifference curves at constant "a" are increasing (absolute value decreasing) as we reduce "p," i.e., that "p" is more highly valued with respect to "a," the less of "p" there is. Another way of looking at this is to note that this implies that "a"'s self-interests are not inferior goods.[4] We could, of course, assume that they are inferior goods (e.g., the manager who insists on a certain level of discretionary return of his perrequisites or emoluments, let us say, as the store of discretionary resources shrinks, no matter what the resulting possible return to the owner's goal of profit), since such cases are not uncommon. For illustrative purposes, however, we assume the normal goods case.

Assume that the agent administers some quantity of these resources, subject to his discretion, which in continuously varying amounts may be either devoted to return to the principal's goal, diverted to the agent's self-goal, or "lost." The agency relation involves the supply or conversion of this resource to increase the return to the principal's goal. Assume that only the agent can convert this resource or otherwise supply it so that it increases return to the principal's goal. The principal will expect some minimal level of performance from his agent, that is, minimal level of resources devoted to return to his goal, below which the principal will seek to discharge his agent. Above this level, $p_i$, the agent has discretion. Here "discretion" means that the agent may choose to perform acts that affect return to the principal's goal. The agent has a similar minimum level, $a_i$, below which he would

[3] For arguments on the comparative behavior of these four types of agents, see Mitnick (1974, 1975).

[4] See Williamson (1964, p. 47). The model of the policing process discussed here was suggested by Williamson's treatment of managerial discretion and Alchian's comments on Williamson. See Alchian (1971). The approach in no way depends on the theory of the firm, however, or on its origin in economics except as the assumptions are similar.

quit as agent. The agent's discretionary choice of returns to his self- and to his other-goal are constrained by a feasibility set of possible combinations of agent-principal receipts, the fiduciary function, p(a). Although we represent it as a function and as a function only of "a," the true fiduciary set may consist of the region below p(a) and depend on other factors, here controlled, besides "a." The assumption of greed, however, makes only the outer boundary of this region relevant, and this is what is indicated by p(a). The shape of the fiduciary function may depend in an organization, for example, on organizational rules, technology, and so on. We assume that through the range of study the agent's preferences will be independent of the fiduciary function. This may of course not be true in general; what is preferred is frequently determined by what one can get. Where the slope of the fiduciary function is positive, an increase in return to the agent's self-goal also implies an increase in return to the principal; where this slope is negative, increase in return to the agent's self-goal means a decrease in return to the principal.



**Figure 1**

Consider the fiduciary function depicted in Figure 1. If the agent is a pure fiduciary, he will select the greatest feasible return to the principal's goal, and will be indifferent regarding the resultant return to any self-goals that may be affected by the agency relation. His indifference curves will be horizontal lines rather than negatively sloped, and he will choose the level of return $p_1$ to the principal's goal. He will be indifferent between the returns to his self-goal at points A and B and between. If, however, the agent is a lexicographical fiduciary, he will select from among those points that equivalently maximize return to the principal that point that yields the highest return to his self-goals, i.e., point B. For the general case where the agent has genuine preferences for both a self- and an other-goal that are not constrained in choice by a fiduciary or other norm, the agent's indifference curves are negatively sloped and determine an "operating point" at the "highest"

10

point of tangency between the indifference curves and the fiduciary function, as in Figure 2 (point A). Note that in Figure 2 the tangency occurs below the point of maximum return to the principal (point B). (See Williamson 1964; Migue and Belanger 1974.) This operating point is allowable only if we assume that the cost to the principal of policing his agent so that he selects the point of maximum principal return exceeds the value of the principal of the additional return that the policing could obtain. (See Alchian 1971; Alchian and Demsetz 1972.)



**Figure 2**

For the purpose of developing a simple model here, assume that $p(a)$ is linear and downward-sloping, i.e., that constant additions to return to the self-goal mean constant reductions in return to the other-goal of the agent. This is similar to Williamson's emoluments model (1974, pp. 49-52), where we will consider lump sum withdrawal of resources rather than a proportional tax as the first policing stage. Assume that the principal for some reason wishes to police his agent. Note that one form of policing would be encouraging the agent to hold the fiduciary norm; we will not, however, consider here the actual mode of policing. Assume that the principal has no outside source of resources to devote to a policing mechanism. He must then divert some of the total discretionary resources potentially available to be distributed towards his own and the agent's ends into a policing apparatus. Note that we have assumed that only the agent may supply or convert the resource so that it increases return to the principal. Alternatively, at this first stage, we may simply consider the effect of withdrawing some part of the discretionary resources without assuming the creation of a policing apparatus. Any withdrawal of resources, or conversion of discretionary to nondiscretionary resources, whether or not policing is intended, is relevant here. Thus, after resource withdrawal but before inauguration of the policing unit — or before that unit begins to have any

effect — we have the situation depicted in Figure 3. For convenience in notation, let $p\text{-}p_i$ be represented by "p," and $a\text{-}a_i$ by "a."



## Figure 3

The policing cost in resources is $p_m\text{-}p_t$. The operating point before is $(a_{mo}, p_{mo})$ and after, $(a_{to}, p_{to})$. $a_{to} < a_{mo}$ because the assumption that "a" is not an inferior good. (In the economic formalism, something similar is called the "income effect.") The fiduciary function is $p = -\dfrac{p_m}{a_m} a + p_m$. Substituting values before and after resource withdrawal, and then subtracting:

$$p_{mo} = -\frac{p_m}{a_m} a_{mo} + p_m$$

$$p_{to} = -\frac{p_m}{a_m} a_{to} + p_t$$

$$(p_{mo} - p_{to}) = -\frac{p_m}{a_m}(a_{mo} - a_{to}) + (p_m - p_t)$$

$$(p_m - p_t) - (p_{mo} - p_{to}) = \frac{p_m}{a_m}(a_{mo} - a_{to}) \tag{1}$$

12

or $$\frac{(p_m - p_t)}{p_m} - \frac{(p_{mo} - p_{to})}{p_m} = \frac{(a_{mo} - a_{to})}{a_m} \tag{2}$$

Eq. (1) suggests the following general law for fiduciary functions of this type: The net gain in agent fidelity is directly proportional to the difference between the discretionary resources withdrawn and the resultant difference between the principal's operating points. Eq. (2) states this in terms of proportionate costs and losses relative to total potential discretionary return to each party. Furthermore, given the assumed preference structure of the agent, if $\frac{p_m}{a_m}$, the absolute value of the slope of the fiduciary function, is $\geqslant 1$ (and it may be some undetermined amount less), that is, if the potential resources that can be realized by the principal exceeds that which could be realized by the agent, the resources withdrawn to be devoted to policing costs (if that is where they go) will always exceed the immediate gain in agent fidelity as measured in resource loss to the agent. For $\frac{p_m}{a_m}$ sufficiently small (that is, linear functions sufficiently more horizontal), however, the resource loss to the agent may exceed the potential resource cost to the principal. Of course, we have not indicated (nor do we know) how the agent and principal will value their respective loss and cost. It is clear that in the limit of some sufficiently large resource withdrawal the principal may obtain a perfectly fiduciary agent, given "room" above the minimum to operate. We note that the theorem above may be subject to empirical study.

According to the above assumptions, a principal may under some specified conditions exchange discretionary resources for an agent acting with greater fidelity to the principal's goal. That is, $(a_{mo}, p_{mo}) > (a_{to}, p_{to})$. This is a paradoxical situation in the short run, since the principal has voluntarily yielded valued resources, and we may ask under what conditions would a *rational* principal do this? Clearly, if the principal has another goal that is satisfied by having a truer agent, he may prefer to do this, i.e., if he has a goal of agent fidelity. The principal may also do this if he expects a net gain as the result of expected future receipts. Only an observer's lack of information about the fact he is observing a stage in a process, and about the principal's full goal set, makes the situation appear paradoxical, of course. Thus we consider now stage two of the policing process: the policing apparatus is in operation; the fiduciary function is altered as the policing unit changes the permissible combinations of resources to the agent and resources to the principal, that is, changes the operating distributive rules. The situation is indicated for the above restrictive model in Figure 4.

The policing device alters the slope of the fiduciary function; the greater the success of the device, the more steeply negative that slope becomes. That is, the policing device makes "a" relatively more expensive to the agent with respect to "p".

## Figure 4

Thus, given the same overall level of "satisfaction," i.e., the same indifference curve, "p" is substituted for "a". (In the economic formalism, something similar is called the "substitution effect.")[5] In Figure 4, we see that the fidelity of the agent is further increased (loss of $a_{to}$ - $a_{so}$ resources to the agent), and the principal has a gain of $p_{so}$ - $p_{to}$ over stage one. The net effect over stages one and two (policing device funded out of total discretionary resources and in operation) is that the principal exchanges operating point $(a_{mo}, p_{mo})$ for point $(a_{so}, p_{so})$. Now although the agent's receipts will always be less than at the start under the present model, the receipts by the principal may be either less than or greater than (or, of course, equal to) receipts at the start. Which of these results obtains depends on the efficiency of the policing device—as indicated by the slope of the fiduciary function — and the shape of the indifference curves of the agent. If $p_{so} < p_{mo}$, we have the policing "paradox" noted above. Again, the principal may be coming out ahead because he also has a goal of agent fidelity, or, again, the principal may have a future reward in mind, which we will consider under stage three. If $p_{so} > p_{mo}$, the principal has

[5]Note that the source of this substitution effect is different from those of the effects discussed in Williamson (1964) and Migue and Belanger (1974), involving as it does changes in the fiduciary function due to organizational changes resulting from policing.

succeeded in taxing the agent not only for paying the total cost of policing, but also into contributing to the principal's receipts.

Now, $a_{so} - a_{to}$ represents the agent's response to policing. Given that the agent is rational, he will clearly yield resources only until the point at which it is no longer rational to do so. Thus the expected value of the negative sanctions of the policing device (including net loss in return to the agent's self-goal) must be greater than or equal to the value of the agent of $a_{so} - a_{to}$. The value of $a_{so} - a_{to}$ may at that point be the expected value of the loss to the agent if the sanctions are applied times the probability that the sanctions will in fact be applied. For example, although the value to the corrupt policeman of the loss of career plus criminal penalties may be very high, he may perceive that the probability of catching him is very low. It may thus be rational for him to respond minimally to the policing attempt.

In stage three, the principal has either succeeded in training his agent to value the principal's goal more highly, perhaps through the policing device, or, if there are multiple agents, replaced those who value "p" very low with agents who value it highly. Thus, for example, the few "rotten apples" in the police department — men with indifference curves with slopes everywhere much more negative (more steeply downward sloping than the average) — have been "weeded out." We then have for our simple model the situation in Figure 5.



**Figure 5**

15

The tangency of the new indifference curves $U_a^2$ is further to the left. The new operating point is $(a_{ro}, p_{ro})$, where $a_{ro} < a_{so} < a_{to} < a_{mo}$, and $p_{ro} > p_{so} > p_{to}$, but $p_{ro}$ either $< p_{mo}$ or $\geqslant p_{mo}$. After stage three the principal may in fact have finally realized his delayed net gain; or he may be left in the situation of the policing "paradox." As above, he may, of course, have a net gain because of the value he places on the increased fidelity of his agent. Note again that the rational principal would not have elected to police his agent if he did not expect a net gain from the attempt. Factors such as error due to uncertainty regarding the effectiveness of the policing mechanism in obtaining the results intended and regarding the choices of the agent, and other factors, may intervene. In addition, the principal may have imperfect information regarding the fiduciary function (which may involve the terms of contract between agent and principal) and the preferences of his agent in various situations. Or the bounded rationality (see March and Simon 1958, Chapter 6; Simon 1957) of the principal and agent may limit their ability to calculate the probable outcomes from electing any action. Or factors relating to strategic interaction between the parties may affect the result at any stage. If the agent at stage three has a hostile reaction to policing,[6] his preferences may shift the opposite way, so that the tangency would be further to the right. The principal would then have a net loss between stages two and three. In order to rationally elect policing, the principal must therefore be able to predict any agent reactions of this type.

We note an additional factor that may be relevant if the policing process occurs over some time. The principal may have to discount his expected return at each stage due to the delay in receiving it. If discounts of this type are required, the principal will demand higher net returns than under the short term or static case. This may necessitate a better or more efficient policing mechanism.

Note also the importance of the fiduciary norm under policing: the fiduciary polices himself. This economizes on policing costs, including costs attached to metering the effects of the policing mechanism and reporting this feedback to the principal. (Cf. those cited by Kaufman 1973.) Thus savings (to whoever pays such costs — most likely, the principal) occurs on both policing costs and specification costs — the costs of choosing agent acts.

In conclusion, we have formulated the policing process in three stages: 1) diversion of resources to policing or other uses; 2) implementation of policing mechanism; and 3) agent's reaction to policing. These stages are, of course, very broad and would require substantial elaboration in a reasonably complete theory of agency. A sense of "paradox" may be experienced by an observer of the policing process because of lack of knowledge regarding the fact that 1) the principal may have a hidden or nonobvious goal, and/or 2) the principal may err, having poor information due to the factors noted above (the principal and agent may also engage in strategic interaction with a paradoxical result occurring because each party may possess poor information on the ultimate effects of their joint

---

[6]On reaction, see, e.g., Brehm (1966), Day and Hamblin (1964), Williamson (1973).

behaviors), and/or 3) the principal may be in the midst of a policing process and expect future returns, the mechanism and dynamics of which are not obvious to the observer.

### III. The Policing "Paradox" and Regulatory Behavior

Under conditions of "organizational slack" (see Cyert and March 1963), the managers of an organization may be viewed as possessing some fixed quantity of discretionary resources to dispose of. (See Williamson 1964). In the firm, for example, revenue that is discretionary may be devoted to increments of profit to the "guardian" stockholders or to increased emoluments or staff that benefit the managers. Similar arguments can be developed for some administrative agencies (cf. Migue and Belanger 1974; Niskanen 1971; Weatherby 1971). In particular, the regulatory agency can be modelled in this fashion.

We assume that the discretionary resources of the agency may go either toward extra satisfaction of some public interest criteria[7] through their indicators, e.g., careful examination of rate applications, or toward increasing the commissioner's status, easing his workload, and insuring the likelihood for him of lucrative future employment, where these rewards may be offered largely by the regulated industry, and may be attained by commission activities that favor the industry.[8] The commissioner's preferences reflect a trade-off between choice of return to self-goals, as represented by activities favoring the industry, and return to public interest criteria. We assume, however, that the commissioner will act more like a self-interest agent than a fiduciary. We expect that, through some range, return to public interest criteria and to self-goals will be instrumental to each other; past this range, self-goals and public interest criteria will be in antinomy. The commissioner enjoys some discretionary return for discretionary return to the public interest, and vice versa; he does gain some added status, for example, for added efforts in behalf of the public interest. But his return from activities generally favoring industry (and against the public interest) soon dominates his sources of return to self-goals. Note that, in the early range, some extra activities favoring the industry may also be in the public interest. The fiduciary function is therefore shaped something like that in Figure 2.

[7]These may be viewed as a set of criteria, aggregated by some calculus into a single dimension. This is obviously difficult to do, if only because of problems of measurement and of operationalization. The welfare tradeoffs involved could require better definition of social values than is now available. In principal, however, such tradeoffs could be established and such a dimension constructed; decisions, after all, are made on such matters all the time. But because of such analytical problems, the subsequent discussion should be thought of as a "gedanken experiment." Employment of the arguments we offer here as explanation does not require that we actually construct operational indices of the variables considered, but only acceptance of the assumption that this could in principle be done.

[8]For arguments regarding the goals of regulators, see e.g., De Alessi (1974), Eckert (1973), Mitnick and Weiss (1974), Noll (1971).

The fixed supply of discretionary resources may include such factors as the commissioner's supply of incentives to distribute within his agency to produce behavior toward one of these ends, i.e., reflect the Barnard (1938) "economy of incentives" model, and the commissioner's disposition of that portion of his own time not already appropriated by certain commission activities. Public interest activist groups may successfully divert part of the discretionary overhead into extensive litigation.[9] If the model developed above is applicable, a policing "paradox" may occur. Because of diversion of resources alone, the first stage of the policing model, we may observe the short run "paradox": The commissioners, in re-distributing their reduced supply of discretionary resources, may act with greater fidelity to the given public interest criteria, i.e., devote less return to industry goals that are instrumental to commissioner self-goals. But the return to the public interest, involving, for example, extensive delays in litigation, may be less than previously. The agency may even be effectively paralyzed, unable to dispose of the cases before it.

Once the diverted resources have been converted into a functioning policing mechanism, perhaps through the adoption of agency regulations requiring use of agency resources to consider public interest questions, the return to that interest may increase. Since the discretionary resources are limited in amount, this involves re-distribution from industry ends instrumental to commissioner goals to public interest ends. It is possible, as we noted in the section above on the policing "paradox," that the increased return to the public interest criteria may still be below its original level. The added delay in litigation may still offset, for example, the gain in respect for public interest criteria reflected in agency decisions.

In time, the commissioners may be "educated" to respect and prefer the public interest group position more, and simultaneously to prefer the goals whose return is mediated by the regulated industry less. The policing mechanism constitutes a shift in the incentive system that may thus, over time, produce

---

[9]Migue and Belanger (1974, p. 46) note that "politicians and through them other parties in the decision-making process can reduce the bureaucrats' margin of discretion." This may occur through constraining their budget. But we note that constraints on true *discretionary* budget may have sources other than constraints on total budget.

Migue and Belanger assume an instantaneous response to an increase in demand: the budget line rises. They assume that demand leads to a rise in discretionary budget. But it is perhaps more likely that in the short run discretionary budget cannot be increased in response to demand; there is a lag due to the political process of approving a new budget. In the short run it is possible, then, that increased demand will simply lead to reduction in the discretionary budget as more of the total budget is consigned to satisfy sharply increased demand. (This could conceivably be represented in the Migue and Belanger model by an increase in minimum "cost.") The fall of the budget line under these conditions may be similar in effect to the policing effect we describe.

Note that extensive litigation may involve both increased cost for the bureau (e.g., more rules to be made and followed) and increased demand (e.g., more rules and decisions are demanded as caseload increases), which would reduce discretionary budget according to our argument above. In either case, according to this modification of the Migue and Belanger model, the budget line would fall, giving rise to policing effects.

genuine changes in preferences. The commissioners may follow a developmental pattern, changing from self-interest agents to lexical self-interest agents to lexical fiduciaries and finally, if unlikely, into pure fiduciaries. The commissioners may also discover that the new incentive system permits them to satisfy their old preferences for self-goals through other means, e.g., through the public interest groups. This would allow them to remain essentially pure self-interest agents (through lexical self-interest could occur 'if choosing in this way were consistently instrumental to self-goals). They may find, for example, that status accrues to the position of noble defender of the public interest, and that this reputation could be converted into remunerative future employment. Because of turnover in commissioners, the new incentive system may tend to "recruit" commissioners whose preferences are satisfied by the new system better than those of the old commissioners. At any rate, in the long run the "paradox" may resolve itself.

Because the "paradox" may involve severe degradation in the actual, net return to the public interest, a question is in order about the rationality of public interest groups that challenge the commissions. The rational principal attempts policing only if he expects a net return, discounted for the intervening adjustment time when his return is below the original level. We suggest that because of bounded rationality and information costs, not to speak of fundamental problems in valuation, the return to the given public interest criteria may be more difficult to measure and predict than certain obvious indicators wrongly assumed to be correlated with return to that interest. (Cf. Niskanen's comments, Migue and Belanger 1974, p. 43.) And these indicators are then relied on in actions that public interest groups consider very rational indeed.[10]

A major indicator is the degree of apparent commissioner fidelity, i.e., how much of discretionary agency resources he diverts to his self-goals, perhaps through the instrumental means of appearing in his conduct to favor the regulated industry. Does the commissioner-agent appear to be a self-interest agent or fiduciary, with gradations between? Note that the real measure of agent fidelity is the difference between what he appropriates for his self-goals and the level of discretionary return to agent self-goals at which discretionary return to the principal is maximized (see points A, B in Figure 2). Since the latter may be unknown, the focus inevitably shifts to simply reducing discretionary return to the agent. If by some chance the organization is already operated to give maximum discretionary return to the principal, then reduction in agent return will result in a reduction in principal return irrespective of any policing "paradox." Public interest groups that object to this "necessary" level of discretionary agent return, e.g., comfortable offices, will thus be acting rather directly against their interest. If, given the information that this is so, such groups persist in such behavior, we would assume either that they are irrational, which is not likely, though possible, or that what they really prefer is

---

[10]Alchian and Demsetz (1972) note in their model of the firm that the difficulty of monitoring marginal productivity leads to the metering of inputs. This result is clearly related to our own and Niskanen's (in "Comment" on Migue and Belanger 1974) arguments.

the appearance of agent fidelity and not the maximization of return to the public interest.

In addition, the legal and other weapons that such groups may use can frequently be directed only against evidence of agency manager fidelity, rather than to insure some given return to the public interest. For example, a group may use litigation to force an agency to prepare an environmental impact statement, or to force an agency that prepared only a pro forma statement to put some content into it. But the group generally cannot litigate to determine that content, i.e., to insure the level of return to the interest they represent that they prefer. Thus even if a group focuses on return to the public interest rather than on fidelity it may only be able to act on commissioner infidelity, hoping the outcome will not be perverse.

At the heart of the policing "paradox" is the observation that agent fidelity in the sense used here does not necessarily correlate with level of principal return. By focusing on apparently improper behavior by the commissioners, such as reliance on the industry for helpful advice, "ex parte" contacts with he industry that may involve exchange of inside information, field visits and spending time with industry representatives rather than with commission staff or "impartial" experts, allocating discretionary resources within the agency to engineering sections with close ties to the industry rather than, say, environmental evaluation sections, and so on, the public interest groups may succeed only in securing an honest agency that doesn't regulate, and, ultimately, a righteous government that cannot govern.

### IV. Conclusion

In conclusion, we have developed a model of policing in the context of agency relations and managerial discretion. The model had three stages: 1) diversion of resources to policing or other uses; 2) implementation of policing mechanism; and 3) agent's reaction to policing. We then applied the model to the case of regulatory behavior. We argued in part that public interest groups are constrained (and perhaps in some cases may elect) to police the manifestations of agent fidelity in the regulatory agencies rather than adherence to public interest criteria. This has possibly paradoxical consequences in that return to public interest criteria may thereby be reduced.

REFERENCES

Alchian, Armen A. "The Basis of Some Recent Advances in the Theory of Management of the Firm." Reprinted on pp. 131-139 in W. Breit and H. M. Hochman (eds.), *Readings in Microeconomics*, Second Edition. New York: Holt, Rinehart and Winston, 1971.

Alchian, Armen A. and Harold Demsetz. "Production, Information Costs, and Economic Organization." *Am. Econ. Rev.*, 62 (December 1972), 777-795.

Barnard, Chester I. *The Functions of the Executive*. Cambridge, Massachusetts: Harvard University Press, 1938, 1968.

Blake, Judith and Kingsley Davis, "Norms, Values, and Sanctions." Pages 456-484 in Robert E. L. Faris (ed.), *Handbook of Modern Sociology*. Chicago: Rand McNally, 1964.

Brehm, Jack. *A Theory of Psychological Reactance*. New York: Academic, 1966.

Curry, R. L., Jr. and L. L. Wade. *A Theory of Political Exchange*. Englewood Cliffs, New Jersey: Prentice-Hall, 1968.

Cyert, Richard M. and James G. March. *A Behavioral Theory of the Firm*. Englewood Cliffs, New Jersey: Prentice-Hall; 1963.

Day, R. C. and R. L. Hamblin. "Some Effects of Close and Punitive Styles of Supervision." *Am. Jour. of Sociol.*, 69 (March 1964), 499-510.

De Alessi, Louis. "An Economic Analysis of Government Ownership and Regulation: Theory and the Evidence from the Electric Power Industry." *Public Choice*, 19 (Fall 1974), 1-42.

Eckert, Ross D. "On the Incentives of Regulators: The Case of Taxicabs." *Public Choice*, 14 (Spring 1973), 83-99.

Goldberg, Victor P. "Consumer Choice, Imperfect Information, and Public Policy." IGA Research Report No. 26 (Davis, California: University of California, Institute of Governmental Affairs, August 1973).

_____. "Regulation and Administered Contracts." Unpublished paper (October, December 1974).

Kaufman, Herbert. *Administrative Feedback*. Washington, D.C.: The Brookings Institution, 1973.

March, James G. and Herbert A. Simon. *Organizations*. New York: Wiley, 1958.

Migue, Jean-Luc and Gerard Belanger. "Toward a General Theory of Managerial Discretion." *Public Choice*, 17 (Spring 1974), 27-47, including "Comment" by William A. Niskanen and "Reply" by Migue and Belanger.

Mitnick, Barry M. "Fiduciary Rationality and Public Policy: The Theory of Agency and Some Consequences." Paper delivered at the 1973 Annual Meeting of the American Political Science Association, New Orleans, Louisiana.

_____. *The Theory of Agency: The Concept of Fiduciary Rationality and Some Consequences*. Unpublished doctoral dissertation, Department of Political Science, University of Pennsylvania (1974).

_____. "The Theory of Agency: The Fiduciary Norm." Paper prepared

for the 1975 Annual Meeting of the American Sociological Association, San Francisco, California.

Mitnick, Barry M. and Charles Weiss, Jr. "The Siting Impasse and a Rational Choice Model of Regulatory Behavior: An Agency for Power Plant Siting." *Jour. Environ. Econ. Mgt.*, 1 (1974), 150-171.

Niskanen, William A., Jr. *Bureaucracy and Representative Government.* Chicago: Aldine-Atherton, 1971.

Noll, Roger G. *Reforming Regulation.* Washington, D.C.: The Brookings Institution, 1971.

Pitkin, Hanna Fenichel. *The Concept of Representation.* Berkeley: University of California Press, 1967.

Riker, William H. *The Theory of Political Coalitions.* New Haven: Yale University Press, 1962.

Ross, Stephen A. "The Economic Theory of Agency: The Principal's Problem." *Am. Econ. Rev.*, 62 (May 1973), 134-139.

_____ . "On the Economic Theory of Agency: The Principle of Similarity." In *Proceedings of the NBER-NSF Conference on Decision Making and Uncertainty*, forthcoming.

Seavey, Warren A. *Handbook of the Law of Agency.* St. Paul, Minnesota: West, 1964.

Simon, Herbert A. "A Behavioral Model of Rational Choice." In Simon, *Models of . Man: Social and Rational.* New York: Wiley, 1957.

Weatherby, James L., Jr. "A Note on Administrative Behavior and Public Policy." *Public Choice*, 11 (Fall 1971), 107-110.

Williamson, Oliver E. *The Economics of Discretionary Behavior: Managerial Objectives in a Theory of the Firm.* Englewood Cliffs, New Jersey: Prentice-Hall, 1964.

_____ . "Some Notes on the Economics of Atmosphere." Fels Discussion Paper No. 29 (Philadelphia: Fels Center of Government, University of Pennsylvania, March 1973).

# THEORY OF THE FIRM: MANAGERIAL BEHAVIOR, AGENCY COSTS AND OWNERSHIP STRUCTURE

## Michael C. JENSEN and William H. MECKLING*

*University of Rochester, Rochester, NY 14627, U.S.A.*

This paper integrates elements from the theory of agency, the theory of property rights and the theory of finance to develop a theory of the ownership structure of the firm. We define the concept of agency costs, show its relationship to the 'separation and control' issue, investigate the nature of the agency costs generated by the existence of debt and outside equity, demonstrate who bears these costs and why, and investigate the Pareto optimality of their existence. We also provide a new definition of the firm, and show how our analysis of the factors influencing the creation and issuance of debt and equity claims is a special case of the supply side of the completeness of markets problem.

> The directors of such [joint-stock] companies, however, being the managers rather of other people's money than of their own, it cannot well be expected, that they should watch over it with the same anxious vigilance with which the partners in a private copartnery frequently watch over their own. Like the stewards of a rich man, they are apt to consider attention to small matters as not for their master's honour, and very easily give themselves a dispensation from having it. Negligence and profusion, therefore, must always prevail, more or less, in the management of the affairs of such a company.
>
> Adam Smith, *The Wealth of Nations*, 1776, Cannan Edition
> (Modern Library, New York, 1937) p. 700.

## 1. Introduction and summary

### 1.1. Motivation of the paper

In this paper we draw on recent progress in the theory of (1) property rights, (2) agency, and (3) finance to develop a theory of ownership structure[1] for the

[1]We do not use the term 'capital structure' because that term usually denotes the relative quantities of bonds, equity, warrants, trade credit, etc., which represent the liabilities of a firm. Our theory implies there is another important dimension to this problem – namely the relative amounts of ownership claims held by insiders (management) and outsiders (investors with no direct role in the management of the firm).

firm. In addition to tying together elements of the theory of each of these three areas, our analysis casts new light on and has implications for a variety of issues in the professional and popular literature such as the definition of the firm, the "separation of ownership and control", the "social responsibility" of business, the definition of a "corporate objective function", the determination of an optimal capital structure, the specification of the content of credit agreements, the theory of organizations, and the supply side of the completeness of markets problem.

Our theory helps explain:

(1) why an entrepreneur or manager in a firm which has a mixed financial structure (containing both debt and outside equity claims) will choose a set of activities for the firm such that the total value of the firm is *less* than it would be if he were the sole owner and why this result is independent of whether the firm operates in monopolistic or competitive product or factor markets;

(2) why his failure to maximize the value of the firm is perfectly consistent with efficiency;

(3) why the sale of common stock is a viable source of capital even though managers do not literally maximize the value of the firm;

(4) why debt was relied upon as a source of capital before debt financing offered any tax advantage relative to equity;

(5) why preferred stock would be issued;

(6) why accounting reports would be provided voluntarily to creditors and stockholders, and why independent auditors would be engaged by management to testify to the accuracy and correctness of such reports;

(7) why lenders often place restrictions on the activities of firms to whom they lend, and why firms would themselves be led to suggest the imposition of such restrictions;

(8) why some industries are characterized by owner-operated firms whose sole outside source of capital is borrowing;

(9) why highly regulated industries such as public utilities or banks will have higher debt equity ratios for equivalent levels of risk than the average non-regulated firm;

(10) why security analysis can be socially productive even if it does not increase portfolio returns to investors.

*1.2. Theory of the firm: An empty box?*

While the literature of economics is replete with references to the "theory of the firm", the material generally subsumed under that heading is not a theory of the firm but actually a theory of markets in which firms are important actors. The firm is a "black box" operated so as to meet the relevant marginal conditions

with respect to inputs and outputs, thereby maximizing profits, or more accurately, present value. Except for a few recent and tentative steps, however, we have no theory which explains how the conflicting objectives of the individual participants are brought into equilibrium so as to yield this result. The limitations of this black box view of the firm have been cited by Adam Smith and Alfred Marshall, among others. More recently, popular and professional debates over the "social responsibility" of corporations, the separation of ownership and control, and the rash of reviews of the literature on the "theory of the firm" have evidenced continuing concern with these issues.[2]

A number of major attempts have been made during recent years to construct a theory of the firm by substituting other models for profit or value maximization; each attempt motivated by a conviction that the latter is inadequate to explain managerial behavior in large corporations.[3] Some of these reformulation attempts have rejected the fundamental principle of maximizing behavior as well as rejecting the more specific profit maximizing model. We retain the notion of maximizing behavior on the part of all individuals in the analysis to follow.[4]

## 1.3. Property rights

An independent stream of research with important implications for the theory of the firm has been stimulated by the pioneering work of Coase, and extended by Alchian, Demsetz and others.[5] A comprehensive survey of this literature is given by Furubotn and Pejovich (1972). While the focus of this research has been "property rights",[6] the subject matter encompassed is far broader than that term suggests. What is important for the problems addressed here is that specification of individual rights determines how costs and rewards will be

[2]Reviews of this literature are given by Peterson (1965), Alchian (1965, 1968), Machlup (1967), Shubik (1970), Cyert and Hedrick (1972), Branch (1973), Preston (1975).

[3]See Williamson (1964, 1970, 1975), Marris (1964), Baumol (1959), Penrose (1958), and Cyert and March (1963). Thorough reviews of these and other contributions are given by Machlup (1961) and Alchian (1965).

Simon (1955) developed a model of human choice incorporating information (search) and computational costs which also has important implications for the behavior of managers. Unfortunately, Simon's work has often been misinterpreted as a denial of maximizing behavior, and misused, especially in the marketing and behavioral science literature. His later use of the term 'satisficing' [Simon (1959)] has undoubtedly contributed to this confusion because it suggests rejection of maximizing behavior rather than maximization subject to costs of information and of decision making.

[4]See Meckling (1976) for a discussion of the fundamental importance of the assumption of resourceful, evaluative, maximizing behavior on the part of individuals in the development of theory. Klein (1976) takes an approach similar to the one we embark on in this paper in his review of the theory of the firm and the law.

[5]See Coase (1937, 1959, 1960), Alchian (1965, 1968), Alchian and Kessel (1962), Demsetz (1967), Alchian and Demsetz (1972), Monsen and Downs (1965), Silver and Auster (1969), and McManus (1975).

[6]Property rights are of course human rights, i.e., rights which are possessed by human beings. The introduction of the wholly false distinction between property rights and human rights in many policy discussions is surely one of the all time great semantic flimflams.

allocated among the participants in any organization. Since the specification of rights is generally effected through contracting (implicit as well as explicit), individual behavior in organizations, including the behavior of managers, will depend upon the nature of these contracts. We focus in this paper on the behavioral implications of the property rights specified in the contracts between the owners and managers of the firm.

### 1.4. Agency costs

Many problems associated with the inadequcy of the current theory of the firm can also be viewed as special cases of the theory of agency relationships in which there is a growing literature.[7] This literature has developed independently of the property rights literature even though the problems with which it is concerned are similar; the approaches are in fact highly complementary to each other.

We define an agency relationship as a contract under which one or more persons (the principal(s)) engage another person (the agent) to perform some service on their behalf which involves delegating some decision making authority to the agent. If both parties to the relationship are utility maximizers there is good reason to believe that the agent will not always act in the best interests of the principal. The *principal* can limit divergences from his interest by establishing appropriate incentives for the agent and by incurring monitoring costs designed to limit the aberrant activities of the agent. In addition in some situations it will pay the *agent* to expend resources (bonding costs) to guarantee that he will not take certain actions which would harm the principal or to ensure that the principal will be compensated if he does take such actions. However, it is generally impossible for the principal or the agent at zero cost to ensure that the agent will make optimal decisions from the principal's viewpoint. In most agency relationships the principal and the agent will incur positive monitoring and bonding costs (non-pecuniary as well as pecuniary), and in addition there will be some divergence between the agent's decisions[8] and those decisions which would maximize the welfare of the principal. The dollar equivalent of the reduction in welfare experienced by the principal due to this divergence is also a cost of the agency relationship, and we refer to this latter cost as the "residual loss". We define *agency costs* as the sum of:

(1)  the monitoring expenditures by the principal,[9]
(2)  the bonding expenditures by the agent,
(3)  the residual loss.

[7] Cf. Berhold (1971), Ross (1973, 1974a), Wilson (1968, 1969), and Heckerman (1975).
[8] Given the optimal monitoring and bonding activities by the principal and agent.
[9] As it is used in this paper the term monitoring includes more than just measuring or observing the behavior of the agent. It includes efforts on the part of the principal to 'control' the behavior of the agent through budget restrictions, compensation policies, operating rules etc.

26

Note also that agency costs arise in any situation involving cooperative effort (such as the co-authoring of this paper) by two or more people even though there is no clear cut principal–agent relationship. Viewed in this light it is clear that our definition of agency costs and their importance to the theory of the firm bears a close relationship to the problem of shirking and monitoring of team production which Alchian and Demsetz (1972) raise in their paper on the theory of the firm.

Since the relationship between the stockholders and manager of a corporation fit the definition of a pure agency relationship it should be no surprise to discover that the issues associated with the "separation of ownership and control" in the modern diffuse ownership corporation are intimately associated with the general problem of agency. We show below that an explanation of why and how the agency costs generated by the corporate form are born leads to a theory of the ownership (or capital) structure of the firm.

Before moving on, however, it is worthwhile to point out the generality of the agency problem. The problem of inducing an "agent" to behave as if he were maximizing the "principal's" welfare is quite general. It exists in all organizations and in all cooperative efforts – at every level of management in firms,[10] in universities, in mutual companies, in cooperatives, in governmental authorities and bureaus, in unions, and in relationships normally classified as agency relationships such as are common in the performing arts and the market for real estate. The development of theories to explain the form which agency costs take in each of these situations (where the contractual relations differ significantly), and how and why they are born will lead to a rich theory of organizations which is now lacking in economics and the social sciences generally. We confine our attention in this paper to only a small part of this general problem – the analysis of agency costs generated by the contractual arrangements between the owners and top management of the corporation.

Our approach to the agency problem here differs fundamentally from most of the existing literature. That literature focuses almost exclusively on the normative aspects of the agency relationship; that is how to structure the contractual relation (including compensation incentives) between the principal and agent to provide appropriate incentives for the agent to make choices which will maximize

[10]As we show below the existence of positive monitoring and bonding costs will result in the manager of a corporation possessing control over some resources which he can allocate (within certain constraints) to satisfy his own preferences. However, to the extent that he must obtain the cooperation of others in order to carry out his tasks (such as divisional vice presidents) and to the extent that he cannot control their behavior perfectly and costlessly they will be able to appropriate some of these resources for their own ends. In short, there are agency costs generated at every level of the organization. Unfortunately, the analysis of these more general organizational issues is even more difficult than that of the 'ownership and control' issue because the nature of the contractual obligations and rights of the parties are much more varied and generally not as well specified in explicit contractual arrangements. Nevertheless, they exist and we believe that extensions of our analysis in these directions show promise of producing insights into a viable theory of organization.

the principal's welfare given that uncertainty and imperfect monitoring exist. We focus almost entirely on the positive aspects of the theory. That is, we assume individuals solve these normative problems and given that only stocks and bonds can be issued as claims, we investigate the incentives faced by each of the parties and the elements entering into the determination of the equilibrium contractual form characterizing the relationship between the manager (i.e., agent) of the firm and the outside equity and debt holders (i.e., principals).

### 1.5. Some general comments on the definition of the firm

Ronald Coase (1937) in his seminal paper on "The Nature of the Firm" pointed out that economics had no positive theory to determine the bounds of the firm. He characterized the bounds of the firm as that range of exchanges over which the market system was suppressed and resource allocation was accomplished instead by authority and direction. He focused on the cost of using markets to effect contracts and exchanges and argued that activities would be included within the firm whenever the costs of using markets were greater than the costs of using direct authority. Alchian and Demsetz (1972) object to the notion that activities within the firm are governed by authority, and correctly emphasize the role of contracts as a vehicle for voluntary exchange. They emphasize the role of monitoring in situations in which there is joint input or team production.[11] We sympathize with the importance they attach to monitoring, but we believe the emphasis which Alchian–Demsetz place on joint input production is too narrow and therefore misleading. Contractual relations are the essence of the firm, not only with employees but with suppliers, customers, creditors, etc. The problem of agency costs and monitoring exists for all of these contracts, independent of whether there is joint production in their sense; i.e., joint production can explain only a small fraction of the behavior of individuals associated with a firm. A detailed examination of these issues is left to another paper.

It is important to recognize that most organizations are simply *legal fictions*[12] *which serve as a nexus for a set of contracting relationships among individuals.* This includes firms, non-profit institutions such as universities, hospitals and foundations, mutual organizations such as mutual savings banks and insurance companies and co-operatives, some private clubs, and even governmental bodies such as cities, states and the Federal government, government enterprises such as TVA, the Post Office, transit systems, etc.

---

[11]They define the classical capitalist firm as a contractual organization of inputs in which there is '(a) joint input production, (b) several input owners, (c) one party who is common to all the contracts of the joint inputs, (d) who has rights to renegotiate any input's contract independently of contracts with other input owners, (e) who holds the residual claim, and (f) who has the right to sell his contractual residual status.'

[12]By legal fiction we mean the artificial construct under the law which allows certain organizations to be treated as individuals.

The private corporation or firm is simply one form of *legal fiction which serves as a nexus for contracting relationships and which is also characterized by the existence of divisible residual claims on the assets and cash flows of the organization which can generally be sold without permission of the other contracting individuals.* While this definition of the firm has little substantive content, emphasizing the essential contractual nature of firms and other organizations focuses attention on a crucial set of questions – why particular sets of contractual relations arise for various types of organizations, what the consequences of these contractual relations are, and how they are affected by changes exogenous to the organization. Viewed this way, it makes little or no sense to try to distinguish those things which are "inside" the firm (or any other organization) from those things that are "outside" of it. There is in a very real sense only a multitude of complex relationships (i.e., contracts) between the legal fiction (the firm) and the owners of labor, material and capital inputs and the consumers of output.[13]

Viewing the firm as the nexus of a set of contracting relationships among individuals also serves to make it clear that the personalization of the firm implied by asking questions such as "what should be the objective function of the firm", or "does the firm have a social responsibility" is seriously misleading. *The firm is not an individual.* It is a legal fiction which serves as a focus for a complex process in which the conflicting objectives of individuals (some of whom may "represent" other oganizations) are brought into equilibrium within a framework of contractual relations. In this sense the "behavior" of the firm is like the behavior of a market; i.e., the outcome of a complex equilibrium process. We seldom fall into the trap of characterizing the wheat or stock market as an individual, but we often make this error by thinking about organizations as if they were persons with motivations and intentions.[14]

---

[13]For example, we ordinarily think of a product as leaving the firm at the time it is sold, but implicitly or explicitly such sales generally carry with them continuing contracts between the firm and the buyer. If the product does not perform as expected the buyer often can and does have a right to satisfaction. Explicit evidence that such implicit contracts do exist is the practice we occasionally observe of specific provision that 'all sales are final.'

[14]This view of the firm points up the important role which the legal system and the law play in social organizations, especially, the organization of economic activity. Statutory laws sets bounds on the kinds of contracts into which individuals and organizations may enter without risking criminal prosecution. The police powers of the state are available and used to enforce performance of contracts or to enforce the collection of damages for non-performance. The courts adjudicate conflicts between contracting parties and establish precedents which form the body of common law. All of these government activities affect both the kinds of contracts executed and the extent to which contracting is relied upon. This in turn determines the usefulness, productivity, profitability and viability of various forms of organization. Moreover, new laws as well as court decisions often can and do change the rights of contracting parties ex post, and they can and do serve as a vehicle for redistribution of wealth. An analysis of some of the implications of these facts is contained in Jensen and Meckling (1976) and we shall not pursue them here.

## 1.6. An overview of the paper

We develop the theory in stages. Sections 2 and 4 provide analyses of the agency costs of equity and debt respectively. These form the major foundation of the theory. Section 3 poses some unanswered questions regarding the existence of the corporate form of organization and examines the role of limited liability. Section 5 provides a synthesis of the basic concepts derived in sections 2–4 into a theory of the corporate ownership structure which takes account of the trade-offs available to the entrepreneur–manager between inside and outside equity and debt. Some qualifications and extensions of the analysis are discussed in section 6, and section 7 contains a brief summary and conclusions.

## 2. The agency costs of outside equity

### 2.1. Overview

In this section we analyze the effect of outside equity on agency costs by comparing the behavior of a manager when he owns 100 percent of the residual claims on a firm to his behavior when he sells off a portion of those claims to outsiders. If a wholly owned firm is managed by the owner, he will make operating decisions which maximize his utility. These decisions will involve not only the benefits he derives from pecuniary returns but also the utility generated by various non-pecuniary aspects of his entrepreneurial activities such as the physical appointments of the office, the attractiveness of the secretarial staff, the level of employee discipline, the kind and amount of charitable contributions, personal relations ("love", "respect", etc.) with employees, a larger than optimal computer to play with, purchase of production inputs from friends, etc. The optimum mix (in the absence of taxes) of the various pecuniary and non-pecuniary benefits is achieved when the marginal utility derived from an additional dollar of expenditure (measured net of any productive effects) is equal for each non-pecuniary item and equal to the marginal utility derived from an additional dollar of after tax purchasing power (wealth).

If the owner–manager sells equity claims on the corporation which are identical to his (i.e., share proportionately in the profits of the firm and have limited liability) agency costs will be generated by the divergence between his interest and those of the outside shareholders, since he will then bear only a fraction of the costs of any non-pecuniary benefits he takes out in maximizing his own utility. If the manager owns only 95 percent of the stock, he will expend resources to the point where the marginal utility derived from a dollar's expenditure of the firm's resources on such items equals the marginal utility of an additional 95 cents in general purchasing power (i.e., *his* share of the wealth reduction) and not one dollar. Such activities, on his part, can be limited (but probably not eliminated) by the expenditure of resources on monitoring activities by the out-

side stockholders. But as we show below, the owner will bear the entire wealth effects of these expected costs so long as the equity market anticipates these effects. Prospective minority shareholders will realize that the owner–manager's interests will diverge somewhat from theirs, hence the price which they will pay for shares will reflect the monitoring costs and the effect of the divergence between the manager's interest and theirs. Nevertheless, ignoring for the moment the possibility of borrowing against his wealth, the owner will find it desirable to bear these costs as long as the welfare increment he experiences from converting his claims on the firm into general purchasing power[15] is large enough to offset them.

As the owner–manager's fraction of the equity falls, his fractional claim on the outcomes falls and this will tend to encourage him to appropriate larger amounts of the corporate resources in the form of perquisites. This also makes it desirable for the minority shareholders to expend more resources in monitoring his behavior. Thus, the wealth costs to the owner of obtaining additional cash in the equity markets rise as his fractional ownership falls.

We shall continue to characterize the agency conflict between the owner–manager and outside shareholders as deriving from the manager's tendency to appropriate perquisites out of the firm's resources for his own consumption. However, we do not mean to leave the impression that this is the only or even the most important source of conflict. Indeed, it is likely that the most important conflict arises from the fact that as the manager's ownership claim falls, his incentive to devote significant effort to creative activities such as searching out new profitable ventures falls. He may in fact avoid such ventures simply because it requires too much trouble or effort on his part to manage or to learn about new technologies. Avoidance of these personal costs and the anxieties that go with them also represent a source of on the job utility to him and it can result in the value of the firm being substantially lower than it otherwise could be.

### 2.2. A simple formal analysis of the sources of agency costs of equity and who bears them

In order to develop some structure for the analysis to follow we make two sets of assumptions. The first set (permanent assumptions) are those which shall carry through almost all of the analysis in sections 2–5. The effects of relaxing some of these are discussed in section 6. The second set (temporary assumptions) are made only for expositional purposes and are relaxed as soon as the basic points have been clarified.

---

[15]For use in consumption, for the diversification of his wealth, or more importantly, for the financing of 'profitable' projects which he could not otherwise finance out of his personal wealth. We deal with these issues below after having developed some of the elementary analytical tools necessary to their solution.

*Permanent assumptions*

(P.1)   All taxes are zero.

(P.2)   No trade credit is available.

(P.3)   All outside equity shares are non-voting.

(P.4)   No complex financial claims such as convertible bonds or preferred stock or warrants can be issued.

(P.5)   No outside owner gains utility from ownership in a firm in any way other than through its effect on his wealth or cash flows.

(P.6)   All dynamic aspects of the multiperiod nature of the problem are ignored by assuming there is only one production–financing decision to be made by the entrepreneur.

(P.7)   The entrepreneur–manager's money wages are held constant throughout the analysis.

(P.8)   There exists a single manager (the peak coordinator) with ownership interest in the firm.

*Temporary assumptions*

(T.1)   The size of the firm is fixed.

(T.2)   No monitoring or bonding activities are possible.

(T.3)   No debt financing through bonds, preferred stock, or personal borrowing (secured or unsecured) is possible.

(T.4)   All elements of the owner–manager's decision problem involving port-folio considerations induced by the presence of uncertainty and the existence of diversifiable risk are ignored.

Define:

$X$   $= \{x_1, x_2, \ldots, x_n\}$ = vector of quantities of all factors and activities within the firm from which the manager derives non-pecuniary bene-fits;[16] the $x_i$ are defined such that his marginal utility is positive for each of them;

$C(X)$ = total dollar cost of providing any given amount of these items;

$P(X)$ = total dollar value to the firm of the productive benefits of $X$;

$B(X) = P(X) - C(X)$ = net dollar benefit to the firm of $X$ ignoring any effects of $X$ on the equilibrium wage of the manager.

Ignoring the effects of $X$ on the manager's utility and therefore on his equili-brium wage rate, the optimum levels of the factors and activities $X$ are defined by $X^*$ such that

$$\frac{\partial B(X^*)}{\partial X^*} = \frac{\partial P(X^*)}{\partial X^*} - \frac{\partial C(X^*)}{\partial X^*} = 0.$$

[16]Such as office space, air conditioning, thickness of the carpets, friendliness of employee relations, etc.

Thus for any vector $X \geqq X^*$ (i.e., where at least one element of $X$ is greater than its corresponding element of $X^*$), $F \equiv B(X^*) - B(X) > 0$ measures the dollar cost to the firm (net of any productive effects) of providing the increment $X - X^*$ of the factors and activities which generate utility to the manager. We assume henceforth that for any given level of cost to the firm, $F$, the vector of factors and activities on which $F$ is spent are those, $\hat{X}$, which yield the manager maximum utility. Thus $F \equiv B(X^*) - B(\hat{X})$.

We have thus far ignored in our discussion the fact that these expenditures on $X$ occur through time and therefore there are tradeoffs to be made across time as well as between alternative elements of $X$. Furthermore, we have ignored the fact that the future expenditures are likely to involve uncertainty (i.e., they are subject to probability distributions) and therefore some allowance must be made for their riskiness. We resolve both of these issues by defining $C$, $P$, $B$, and $F$ to be the *current market values* of the sequence of probability distributions on the period by period cash flows involved.[17]

Given the definition of $F$ as the current market value of the stream of manager's expenditures on non-pecuniary benefits we represent the constraint which a single owner–manager faces in deciding how much non-pecuniary income he will extract from the firm by the line $\bar{V}F$ in fig. 1. This is analogous to a budget constraint. The market value of the firm is measured along the vertical axis and the market value of the manager's stream of expenditures on non-pecuniary benefits, $F$, are measured along the horizontal axis. $0\bar{V}$ is the value of the firm when the amount of non-pecuniary income consumed is zero. By definition $\bar{V}$ is the maximum market value of the cash flows generated by the firm for a given money wage for the manager when the manager's consumption of non-pecuniary benefits are zero. At this point all the factors and activities within the firm which generate utility for the manager are at the level $X^*$ defined above. There is a different budget constraint $\bar{V}F$ for each possible scale of the firm (i.e., level of investment, $I$) and for alternative levels of money wage, $W$, for the manager. For the moment we pick an arbitrary level of investment (which we assume has already been made) and hold the scale of the firm constant at this level. We also assume that the manager's money wage is fixed at the level $W^*$ which represents the current market value of his wage contract[18] in the optimal compensation package which consists of both wages, $W^*$, and non-pecuniary benefits, $F^*$. Since one dollar of current value of non-pecuniary benefits withdrawn from the firm by the manager reduces the market value of the firm by \$1, by definition, the slope of $\bar{V}F$ is $-1$.

[17]And again we assume that for any given market value of these costs, $F$, to the firm the allocation across time and across alternative probability distributions is such that the manager's current expected utility is at a maximum.

[18]At this stage when we are considering a 100% owner-managed firm the notion of a 'wage contract' with himself has no content. However, the 100% owner-managed case is only an expositional device used in passing to illustrate a number of points in the analysis, and we ask the reader to bear with us briefly while we lay out the structure for the more interesting partial ownership case where such a contract does have substance.

The owner–manager's tastes for wealth and non-pecuniary benefits is represented in fig. 1 by a system of indifference curves, $U_1$, $U_2$, etc.[19] The indifference curves will be convex as drawn as long as the owner–manager's marginal rate of



MARKET VALUE OF THE STREAM OF MANAGER'S EXPENDITURES
ON NON-PECUNIARY BENEFITS

Fig. 1. The value of the firm ($V$) and the level of non-pecuniary benefits consumed ($F$) when the fraction of outside equity is $(1-\alpha)V$, and $U_j$ $(j = 1, 2, 3)$ represents owner's indifference curves between wealth and non-pecuniary benefits.

substitution between non-pecuniary benefits and wealth diminishes with increasing levels of the benefits. For the 100 percent owner–manager, this presumes that there are not perfect substitutes for these benefits available on the outside, i.e., to some extent they are job specific. For the fractional owner–manager this presumes the benefits cannot be turned into general purchasing power at a constant price.[20]

[19]The manager's utility function is actually defined over wealth and the future time sequence of vectors of quantities of non-pecuniary benefits, $X_t$. Although the setting of his problem is somewhat different, Fama (1970b, 1972) analyzes the conditions under which these preferences can be represented as a derived utility function defined as a function of the money value of the expenditures (in our notation $F$) on these goods conditional on the prices of goods. Such a utility function incorporates the optimization going on in the background which define $\hat{X}$ discussed above for a given $F$. In the more general case where we allow a time series of consumption, $\hat{X}_t$, the optimization is being carried out across both time and the components of $X_t$ for fixed $F$.

[20]This excludes, for instance, (a) the case where the manager is allowed to expend corporate resources on anything he pleases in which case $F$ would be a perfect substitute for wealth, or (b) the case where he can 'steal' cash (or other marketable assets) with constant returns to scale – if he could the indifference curves would be straight lines with slope determined by the fence commission.

When the owner has 100 percent of the equity, the value of the firm will be $V^*$ where indifference curve $U_2$ is tangent to $VF$, and the level of non-pecuniary benefits consumed is $F^*$. If the owner sells the entire equity but remains as manager, and if the equity buyer can, at zero cost, force the old owner (as manager) to take the same level of non-pecuniary benefits as he did as owner, then $V^*$ is the price the new owner will be willing to pay for the entire equity.[21]

In general, however, we would not expect the new owner to be able to enforce identical behavior on the old owner at zero costs. If the old owner sells a fraction of the firm to an outsider, he, as manager, will no longer bear the full cost of any non-pecuniary benefits he consumes. Suppose the owner sells a share of the firm, $1 - \alpha$, $(0 < \alpha < 1)$ and retains for himself a share, $\alpha$. If the prospective buyer believes that the owner–manager will consume the same level of non-pecuniary benefits as he did as full owner, the buyer will be willing to pay $(1 - \alpha)V^*$ for a fraction $(1 - \alpha)$ of the equity. Given that an outsider now holds a claim to $(1 - \alpha)$ of the equity, however, the *cost* to the owner–manager of consuming \$1 of non-pecuniary benefits in the firm will no longer be \$1. Instead, it will be $\alpha \times \$1$. If the prospective buyer actually paid $(1 - \alpha)V^*$ for his share of the equity, and if thereafter the manager could choose whatever level of non-pecuniary benefits he liked, his budget constraint would be $V_1 P_1$ in fig. 1 and has a slope equal to $-\alpha$. Including the payment the owner receives from the buyer as part of the owner's post-sale wealth, his budget constraint, $V_1 P_1$, must pass through $D$, since he can if he wishes have the same wealth and level of non-pecuniary consumption he consumed as full owner.    ·

But if the owner–manager is free to choose the level of perquisites, $F$, subject only to the loss in wealth he incurs as a part owner, his welfare will be maximized by increasing his consumption of non-pecuniary benefits. He will move to point $A$ where $V_1 P_1$ is tangent to $U_1$ representing a higher level of utility. The value of the firm falls from $V^*$, to $V^0$, i.e., by the amount of the cost to the firm of the increased non-pecuniary expenditures, and the owner–manager's consumption of non-pecuniary benefits rises from $F^*$ to $F^0$.

[21] Point $D$ defines the fringe benefits in the optimal pay package since the value to the manager of the fringe benefits $F^*$ is greater than the cost of providing them as is evidenced by the fact that $U_2$ is steeper to the left of $D$ than the budget constraint with slope equal to $-1$.

That $D$ is indeed the optimal pay package can easily be seen in this situation since if the conditions of the sale to a new owner specified that the manager would receive no fringe benefits after the sale he would require a payment equal to $V_3$ to compensate him for the sacrifice of his claims to $V^*$ and fringe benefits amounting to $F^*$ (the latter with total value to him of $V_3 - V^*$). But if $F = 0$, the value of the firm is only $V$. Therefore, if monitoring costs were zero the sale would take place at $V^*$ with provision for a pay package which included fringe benefits of $F^*$ for the manager.

This discussion seems to indicate there are two values for the 'firm', $V_3$ and $V^*$. This is not the case if we realize that $V^*$ is the value of the right to be the residual claimant on the cash flows of the firm and $V_3 - V^*$ is the value of the managerial rights, i.e., the right to make the operating decisions which include access to $F^*$. There is at least one other right which has value which plays no formal role in the analysis as yet – the value of the control right. By control right we mean the right to hire and fire the manager and we leave this issue to a future paper.

If the equity market is characterized by rational expectations the buyers will be aware that the owner will increase his non-pecuniary consumption when his ownership share is reduced. If the owner's response function is known or if the equity market makes unbiased estimates of the owner's response to the changed incentives, the buyer will not pay $(1-\alpha)V^*$ for $(1-\alpha)$ of the equity.

*Theorem.   For a claim on the firm of $(1-\alpha)$ the outsider will pay only $(1-\alpha)$ times the value he expects the firm to have given the induced change in the behavior of the owner–manager.*

*Proof.*   For simplicity we ignore any element of uncertainty introduced by the lack of perfect knowledge of the owner–manager's response function. Such uncertainty will not affect the final solution if the equity market is large as long as the estimates are rational (i.e., unbiased) and the errors are independent across firms. The latter condition assures that this risk is diversifiable and therefore equilibrium prices will equal the expected values.

Let $W$ represent the owner's total wealth after he has sold a claim equal to $1-\alpha$ of the equity to an outsider. $W$ has two components. One is the payment, $S_o$, made by the outsider for $1-\alpha$ of the equity; the rest, $S_i$, is the value of the owner's (i.e., insider's) share of the firm, so that $W$, the owner's wealth, is given by

$$W = S_o + S_i = S_o + \alpha V(F, \alpha),$$

where $V(F, \alpha)$ represents the value of the firm given that the manager's fractional ownership share is $\alpha$ and that he consumes perquisites with current market value of $F$. Let $V_2 P_2$, with a slope of $-\alpha$ represent the tradeoff the owner–manager faces between non-pecuniary benefits and his wealth after the sale. Given that the owner has decided to sell a claim $1-\alpha$ of the firm, his welfare will be maximized when $V_2 P_2$ is tangent to some indifference curve such as $U_3$ in fig. 1. A price for a claim of $(1-\alpha)$ on the firm that is satisfactory to both the buyer and the seller will require that this tangency occur along $\overline{V}F$, i.e., that the value of the firm must be $V'$. To show this, assume that such is not the case – that the tangency occurs to the left of the point $B$ on the line $\overline{V}F$. Then, since the slope of $V_2 P_2$ is negative, the value of the firm will be larger than $V'$. The owner–manager's choice of this lower level of consumption of non-pecuniary benefits will imply a higher value both to the firm as a whole and to the fraction of the firm $(1-\alpha)$ which the outsider has acquired; that is, $(1-\alpha)V' > S_o$. From the owner's viewpoint, he has sold $1-\alpha$ of the firm for less than he could have, given the (assumed) lower level of non-pecuniary benefits he enjoys. On the other hand, if the tangency point $B$ is to the right of the line $\overline{V}F$, the owner–manager's higher consumption of non-pecuniary benefits means the value of the firm is less than $V'$, and hence $(1-\alpha)V(F, \alpha) < S_o = (1-\alpha)V'$. The outside owner then has paid more for his share of the equity than it is worth. $S_o$ will be a mutually satisfactory

36

price if and only if $(1-\alpha)V' = S_o$. But this means that the owner's post-sale wealth is equal to the (reduced) value of the firm $V'$, since

$$W = S_o + \alpha V' = (1-\alpha)V' + \alpha V' = V'.$$

Q.E.D.

The requirement that $V'$ and $F'$ fall on $\overline{V}F$ is thus equivalent to requiring that the value of the claim acquired by the outside buyer be equal to the amount he pays for it and conversely for the owner. *This means that the decline in the total value of the firm ($V^* - V'$) is entirely imposed on the owner–manager.* His total wealth after the sale of $(1-\alpha)$ of the equity is $V'$ and the decline in his wealth is $V^* - V'$.

The distance $V^* - V'$ is the reduction in the market value of the firm engendered by the agency relationship and is a measure of the "residual loss" defined earlier. In this simple example the residual loss represents the total agency costs engendered by the sale of outside equity because monitoring and bonding activities have not been allowed. The welfare loss the owner incurs is less than the residual loss by the value to him of the increase in non-pecuniary benefits ($F' - F^*$). In fig. 1 the difference between the intercepts on the $Y$ axis of the two indifference curves $U_2$ and $U_3$ is a measure of the owner–manager's welfare loss due to the incurrence of agency costs,[22] and he would sell such a claim only if the increment in welfare he achieves by using the cash amounting to $(1-\alpha)V'$ for other things was worth more to him than this amount of wealth.

### 2.3. Determination of the optimal scale of the firm

*The case of all equity financing.* Consider the problem faced by an entrepreneur with initial pecuniary wealth, $W$, and monopoly access to a project requiring investment outlay, $I$, subject to diminishing returns to scale in $I$. Fig. 2 portrays the solution to the optimal scale of the firm taking into account the agency costs associated with the existence of outside equity. The axes are as defined in fig. 1 except we now plot on the vertical axis the total wealth of the owner, i.e., his initial wealth, $W$, plus $V(I) - I$, the net increment in wealth he obtains from exploitation of his investment opportunities. The market value of the firm, $V = V(I, F)$, is now a function of the level of investment, $I$, and the current market value of the manager's expenditures of the firm's resources on non-pecuniary benefits, $F$. Let $\overline{V}(I)$ represent the value of the firm as a function of the level of investment when the manager's expenditures on non-pecuniary benefits, $F$, are zero. The schedule with intercept labeled $W + [\overline{V}(I^*) - I^*)]$ and

---

[22]The distance $V^* - V'$ is a measure of what we will define as the gross agency costs. The distance $V_3 - V_4$ is a measure of what we call net agency costs, and it is this measure of agency costs which will be minimized by the manager in the general case where we allow investment to change.

slope equal to −1 in fig. 2 represents the locus of combinations of post-investment wealth and dollar cost to the firm of non-pecuniary benefits which are available to the manager when investment is carried to the value maximizing



Fig. 2. Determination of the optimal scale of the firm in the case where no monitoring takes place. Point $C$ denotes optimum investment, $I^*$, and non-pecuniary benefits, $F^*$, when investment is 100% financed by entrepreneur. Point $D$ denotes optimum investment, $I'$, and non-pecuniary benefits, $F$, when outside equity financing is used to help finance the investment and the entrepreneur owns a fraction $\alpha'$ of the firm. The distance $A$ measures the gross agency costs.

point, $I^*$. At this point $\Delta V(I) - \Delta I = 0$. If the manager's wealth were large enough to cover the investment required to reach this scale of operation, $I^*$, he would consume $F^*$ in non-pecuniary benefits and have pecuniary wealth with value $W + V^* - I^*$. However, if outside financing is required to cover the investment he will not reach this point if monitoring costs are non-zero.[23]

The expansion path $OZBC$ represents the equilibrium combinations of wealth and non-pecuniary benefits, $F$, which the manager could obtain if he had enough

[23]$I^*$ is the value maximizing and Pareto Optimum investment level which results from the traditional analysis of the corporate investment decision if the firm operates in perfectly competitive capital and product markets and the agency cost problems discussed here are ignored. See Debreu (1959, ch. 7), Jensen and Long (1972), Long (1972), Merton and Subrahmanyam (1974), Hirshleifer (1958, 1970), and Fama and Miller (1972).

personal wealth to finance all levels of investment up to $I^*$. It is the locus of points such as $Z$ and $C$ which represent the equilibrium position for the 100 percent owner–manager at each possible level of investment, $I$. As $I$ increases we move up the expansion path to the point $C$ where $V(I)-I$ is at a maximum. Additional investment beyond this point reduces the net value of the firm, and as it does the equilibrium path of the manager's wealth and non-pecuniary benefits retraces (in the reverse direction) the curve $OZBC$. We draw the path as a smooth concave function only as a matter of convenience.

If the manager obtained outside financing and if there were zero costs to the agency relationship (perhaps because monitoring costs were zero) the expansion path would also be represented by $OZBC$. Therefore, this path represents what we might call the "idealized" solutions, i.e., those which would occur in the absence of agency costs.

Assume the manager has sufficient personal wealth to completely finance the firm only up to investment level $I_1$ which puts him at point $Z$. At this point $W = I_1$. To increase the size of the firm beyond this point he must obtain outside financing to cover the additional investment required, and this means reducing his fractional ownership. When he does this he incurs agency costs, and the lower is his ownership fraction the larger are the agency costs he incurs. However, if the investments requiring outside financing are sufficiently profitable his welfare will continue to increase.

The expansion path $ZEDHL$ in fig. 2 portrays one possible path of the equilibrium levels of the owner's non-pecuniary benefits and wealth at each possible level of investment higher than $I_1$. This path is the locus of points such as $E$ or $D$ where (1) the manager's indifference curve is tangent to a line with slope equal to $-\alpha$ (his fractional claim on the firm at that level of investment), and (2) the tangency occurs on the "budget constraint" with slope $= -1$ for the firm value and non-pecuniary benefit tradeoff at the same level of investment.[24] As we move along $ZEDHL$ his fractional claim on the firm continues

---

[24] Each equilibrium point such as that at $E$ is characterized by $(\hat{a}, \hat{F}, \hat{W}_T)$ where $\hat{W}_T$ is the entrepreneur's post-investment financing wealth. Such an equilibrium must satisfy each of the following four conditions:

(1) $\qquad \hat{W}_T + \hat{F} = V(I) + W - I = V(I) - K,$

where $K \equiv I - W$ is the amount of outside financing required to make the investment $I$. If this condition is not satisfied there is an uncompensated wealth transfer (in one direction or the other) between the entrepreneur and outside equity buyers.

(2) $\qquad U_F(\hat{W}_T, \hat{F})/U_{W_T}(\hat{W}_T, \hat{F}) = \hat{a},$

where $U$ is the entrepreneur's utility function on wealth and perquisites, $U_F$ and $U_{W_T}$ are marginal utilities and $\hat{a}$ is the manager's share of the firm.

(3) $\qquad (1-\hat{a})V(I) = (1-\hat{a})[\hat{V}(I)-\hat{F}] \geqq K,$

which says the funds received from outsiders are at least equal to $K$, the minimum required outside financing.

(4) Among all points $(\hat{a}, \hat{F}, \hat{W}_T)$ satisfying conditions (1)–(3), $(\alpha, \hat{F}, W_T)$ gives the manager highest utility. This implies that $(\hat{a}, \hat{F}, \hat{W}_T)$ satisfy condition (3) as an equality.

to fall as he raises larger amounts of outside capital. This expansion path represents his complete opportunity set for combinations of wealth and non-pecuniary benefits given the existence of the costs of the agency relationship with the outside equity holders. Point $D$, where this opportunity set is tangent to an indifference curve, represents the solution which maximizes his welfare. At this point, the level of investment is $I'$, his fractional ownership share in the firm is $\alpha'$, his wealth is $W + V' - I'$, and he consumes a stream of non-pecuniary benefits with current market value of $F'$. The gross agency costs (denoted by $A$) are equal to $(V^* - I^*) - (V' - I')$. Given that no monitoring is possible, $I'$ is the socially optimal level of investment as well as the privately optimal level.

We can characterize the optimal level of investment as that point, $I'$ which satisfies the following condition for small changes:

$$\Delta V - \Delta I + \alpha' \Delta F = 0. \tag{1}$$

$\Delta V - \Delta I$ is the change in the net market value of the firm, and $\alpha' \Delta F$ is the dollar value to the manager of the incremental fringe benefits he consumes (which cost the firm $\Delta F$ dollars).[25] Furthermore, recognizing that $V = \bar{V} - F$, where $\bar{V}$ is the value of the firm at any level of investment when $F = 0$, we can substitute into the optimum condition to get

$$(\Delta \bar{V} - \Delta I) - (1 - \alpha') \Delta F = 0 \tag{3}$$

as an alternative expression for determining the optimum level of investment.

The idealized or zero agency cost solution, $I^*$, is given by the condition $(\Delta \bar{V} - \Delta I) = 0$, and since $\Delta F$ is positive the actual welfare maximizing level of investment $I'$ will be less than $I^*$, because $(\Delta \bar{V} - \Delta I)$ must be positive at $I'$ if (3) is to be satisfied. Since $-\alpha'$ is the slope of the indifference curve at the optimum and therefore represents the manager's demand price for incremental non-pecuniary benefits, $\Delta F$, we know that $\alpha' \Delta F$ is the dollar value to him of an increment of fringe benefits costing the firm $\Delta F$ dollars. The term $(1 - \alpha') \Delta F$ thus measures the dollar "loss" to the firm (and himself) of an additional $\Delta F$ dollars spent on non-pecuniary benefits. The term $\Delta \bar{V} - \Delta I$ is the gross increment in the value of the firm ignoring any changes in the consumption of non-pecuniary benefits. Thus, the manager stops increasing the size of the firm when the gross

---

[25] *Proof.* Note that the slope of the expansion path (or locus of equilibrium points) at any point is $(\Delta V - \Delta I)/\Delta F$ and at the optimum level of investment this must be equal to the slope of the manager's indifference curve between wealth and market value of fringe benefits, $F$. Furthermore, in the absence of monitoring, the slope of the indifference curve, $\Delta W/\Delta F$, at the equilibrium point, $D$, must be equal to $-\alpha'$. Thus,

$$(\Delta V - \Delta I)/\Delta F = -\alpha' \tag{2}$$

is the condition for the optimal scale of investment and this implies condition (1) holds for small changes at the optimum level of investment, $I'$.

increment in value is just offset by the incremental "loss" involved in the consumption of additional fringe benefits due to his declining fractional interest in the firm.[26]

## 2.4. The role of monitoring and bonding activities in reducing agency costs

In the above analysis we have ignored the potential for controlling the behavior of the owner–manager through monitoring and other control activities. In practice, it is usually possible by expending resources to alter the opportunity the owner–manager has for capturing non-pecuniary benefits. These methods include auditing, formal control systems, budget restrictions, and the establishment of incentive compensation systems which serve to more closely identify the manager's interests with those of the outside equity holders, etc. Fig. 3 portrays the effects of monitoring and other control activities in the simple situation portrayed in fig. 1. Figs. 1 and 3 are identical except for the curve *BCE* in fig. 3 which depicts a "budget constraint" derived when monitoring possibilities are taken into account. Without monitoring, and with outside equity of $(1 - \alpha)$, the value of the firm will be $V'$ and non-pecuniary expenditures $F'$. By incurring monitoring costs, $M$, the equity holders can restrict the manager's consumption of perquisites to amounts less than $F'$. Let $F(M, \alpha)$ denote the maximum perquisites the manager can consume for alternative levels of monitoring expenditures, $M$, given his ownership share $\alpha$. We assume that increases in monitoring reduce $F$, and reduce it at a decreasing rate, i.e., $\partial F/\partial M < 0$ and $\partial^2 F/\partial M^2 > 0$.

Since the current value of expected future monitoring expenditures by the outside equity holders reduce the value of any given claim on the firm to them dollar for dollar, the outside equity holders will take this into account in determining the maximum price they will pay for any given fraction of the firm's

---

[26]Since the manager's indifference curves are negatively sloped we know that the optimum scale of the firm, point *D*, will occur in the region where the expansion path has negative slope, i.e., the market value of the firm will be declining and the *gross* agency costs, *A*, will be increasing and thus, the manager will not minimize them in making the investment decision (even though he will minimize them for any *given* level of investment). However, we define the *net* agency cost as the dollar equivalent of the welfare loss the manager experiences because of the agency relationship evaluated at $F = 0$ (the vertical distance between the intercepts on the *Y* axis of the two indifference curves on which points *C* and *D* lie). The optimum solution, *I'*, does satisfy the condition that net agency costs are minimized. But this simply amounts to a restatement of the assumption that the manager maximizes his welfare.

Finally, it is possible for the solution point *D* to be a corner solution and in this case the value of the firm will not be declining. Such a corner solution can occur, for instance, if the manager's marginal rate of substitution between *F* and wealth falls to zero fast enough as we move up the expansion path, or if the investment projects are 'sufficiently' profitable. In these cases the expansion path will have a corner which lies on the maximum value budget constraint with intercept $V(I^*) - I^*$, and the level of investment will be equal to the idealized optimum, *I\**. However, the market value of the residual claims will be less than $V^*$ because the manager's consumption of perquisites will be larger than $F^*$, the zero agency cost level.

equity. Therefore, given positive monitoring activity the value of the firm is given by $V = \bar{V} - F(M, \alpha) - M$ and the locus of these points for various levels of $M$ and for a given level of $\alpha$ lie on the line $BCE$ in fig. 3. The vertical difference between the $\bar{V}F$ and $BCE$ curves is $M$, the current market value of the future monitoring expenditures.



Fig. 3. The value of the firm $(V)$ and level of non-pecuniary benefits $(F)$ when outside equity is $(1-\alpha)$, $U_1$, $U_2$, $U_3$ represent owner's indifference curves between wealth and non-pecuniary benefits, and monitoring (or bonding) activities impose opportunity set $BCE$ as the tradeoff constraint facing the owner.

If it is possible for the outside equity holders to make these monitoring expenditures and thereby to impose the reductions in the owner–manager's consumption of $F$, he will voluntarily enter into a contract with the outside equity holders which gives them the rights to restrict his consumption of non-pecuniary items to $F''$. He finds this desirable because it will cause the value of the firm to rise to $V''$. Given the contract, the optimal monitoring expenditure on the part of the outsiders, $M$, is the amount $D - C$. The entire increase in the value of the firm that accrues will be reflected in the owner's wealth, but his welfare will be increased by less than this because he forgoes some non-pecuniary benefits he previously enjoyed.

If the equity market is competitive and makes unbiased estimates of the effects

of the monitoring expenditures on $F$ and $V$, potential buyers will be indifferent between the following two contracts:

(i)   Purchase of a share $(1-\alpha)$ of the firm at a total price of $(1-\alpha)V'$ and no rights to monitor or control the manager's consumption of perquisites.

(ii)  Purchase of a share $(1-\alpha)$ of the firm at a total price of $(1-\alpha)V''$ and the right to expend resources up to an amount equal to $D-C$ which will limit the owner–manager's consumption of perquisites to $F;$.

Given contract (ii) the outside shareholders would find it desirable to monitor to the full rights of their contract because it will pay them to do so. However, if the equity market is competitive the total benefits (net of the monitoring costs) will be capitalized into the price of the claims. Thus, not surprisingly, the owner-- manager reaps all the benefits of the opportunity to write and sell the monitoring contract.[27]

*An analysis of bonding expenditures.* We can also see from the analysis of fig. 3 that it makes no difference who actually makes the monitoring expendi- tures – the owner bears the full amount of these costs as a wealth reduction in all cases. Suppose that the owner–manager could expend resources to guarantee to the outside equity holders that he would limit his activities which cost the firm $F$. We call these expenditures "bonding costs", and they would take such forms as contractual guarantees to have the financial accounts audited by a public account, explicit bonding against malfeasance on the part of the manager, and contractual limitations on the manager's decision making power (which impose costs on the firm because they limit his ability to take full advantage of some profitable opportunities as well as limiting his ability to harm the stock- holders while making himself better off).

If the incurrence of the bonding costs were entirely under the control of the manager and if they yielded the same opportunity set *BCE* for him in fig. 3, he would incur them in amount $D-C$. This would limit his consumption of

[27]The careful reader will note that point $C$ will be the equilibrium point only if the contract between the manager and outside equity holders specifies with no ambiguity that they have the right to monitor to limit his consumption of perquisites to an amount no less than $F''$. If any ambiguity regarding these rights exists in this contract then another source of agency costs arises which is symmetrical to our original problem. If they could do so the outside equity holders would monitor to the point where the net value of *their* holdings, $(1-\alpha)V-M$, was maximized, and this would occur when $(\partial V/\partial M)(1-\alpha)-1 = 0$ which would be at some point between points $C$ and $E$ in fig. 3. Point $E$ denotes the point where the value of the firm net of the monitoring costs is at a maximum, i.e. where $\partial V/\partial M-1 = 0$. But the manager would be worse off than in the zero monitoring solution if the point where $(1-\alpha)V-M$ was at a maxi- mum were to the left of the intersection between *BCE* and the indifference curve $U_3$ passing through point $B$ (which denotes the zero monitoring level of welfare). Thus if the manager could not eliminate enough of the ambiguity in the contract to push the equilibrium to the right of the intersection of the curve *BCE* with indifference curve $U_3$ he would not engage in any contract which allowed monitoring.

perquisites to $F''$ from $F'$, and the solution is exactly the same as if the outside equity holders had performed the monitoring. The manager finds it in his interest to incur these costs as long as the net increments in his wealth which they generate (by reducing the agency costs and therefore increasing the value of the firm) are more valuable than the perquisites given up. This optimum occurs at point $C$ in both cases under our assumption that the bonding expenditures yield the same opportunity set as the monitoring expenditures. In general, of course, it will pay the owner–manager to engage in bonding activities and to write contracts which allow monitoring as long as the marginal benefits of each are greater than their marginal cost.

*Optimal scale of the firm in the presence of monitoring and bonding activities.* If we allow the outside owners to engage in (costly) monitoring activities to limit the manager's expenditures on non-pecuniary benefits and allow the manager to engage in bonding activities to guarantee to the outside owners that he will limit his consumption of $F$ we get an expansion path such as that illustrated in fig. 4 on which $Z$ and $G$ lie. We have assumed in drawing fig. 4 that the cost functions involved in monitoring and bonding are such that some positive levels of the activities are desirable, i.e., yield benefits greater than their cost. If this is not true the expansion path generated by the expenditure of resources on these activities would lie below $ZD$ and no such activity would take place at any level of investment. Points $Z$, $C$, and $D$ and the two expansion paths they lie on are identical to those portrayed in fig. 2. Points $Z$ and $C$ lie on the 100 percent ownership expansion path, and points $Z$ and $D$ lie on the fractional ownership, zero monitoring and bonding activity expansion path.

The path on which points $Z$ and $G$ lie is the one given by the locus of equilibrium points for alternative levels of investment characterized by the point labeled $C$ in fig. 3 which denotes the optimal level of monitoring and bonding activity and resulting values of the firm and non-pecuniary benefits to the manager given a fixed level of investment. If any monitoring or bonding is cost effective the expansion path on which $Z$ and $G$ lie must be above the nonmonitoring expansion path over some range. Furthermore, if it lies anywhere to the right of the indifference curve passing through point $D$ (the zero monitoring–bonding solution) the final solution to the problem will involve positive amounts of monitoring and/or bonding activities. Based on the discussion above we know that as long as the contracts between the manager and outsiders are unambiguous regarding the rights of the respective parties the final solution will be at that point where the new expansion path is just tangent to the highest indifference curve. At this point the optimal level of monitoring and bonding expenditures are $M''$ and $b''$; the manager's post-investment-financing wealth is given by $W + V'' - I'' - M'' - b''$ and his non-pecuniary benefits are $F''$. The total gross agency costs, $A$, are given by $A(M'', b'', \alpha'', I'') = (V^* - I^*) - (V'' - I'' - M'' - b'')$.

## 2.5. Pareto optimality and agency costs in manager-operated firms

In general we expect to observe both bonding and external monitoring activities, and the incentives are such that the levels of these activities will satisfy the conditions of efficiency. They will not, however, result in the firm being run in a manner so as to maximize its value. The difference between $V^*$, the efficient solution under zero monitoring and bonding costs (and therefore zero agency



Fig. 4. Determination of optimal scale of the firm allowing for monitoring and bonding activities. Optimal monitoring costs are $M''$ and bonding costs are $b''$ and the equilibrium scale of firm, manager's wealth and consumption of non-pecuniary benefits are at point $G$.

costs), and $V''$, the value of the firm given positive monitoring costs, are the total gross agency costs defined earlier in the introduction. These are the costs of the "separation of ownership and control" which Adam Smith focused on in the passage quoted at the beginning of this paper and which Berle and Means (1932) popularized 157 years later. The solutions outlined above to our highly simplified problem imply that agency costs will be positive as long as monitoring costs are positive – which they certainly are.

The reduced value of the firm caused by the manager's consumption of perquisites outlined above is "non-optimal" or inefficient only in comparison

45

to a world in which we could obtain compliance of the agent to the principal's wishes at zero cost or in comparison to a *hypothetical* world in which the agency costs were lower. But these costs (monitoring and bonding costs and 'residual loss') are an unavoidable result of the agency relationship. Furthermore, since they are borne entirely by the decision maker (in this case the original owner) responsible for creating the relationship he has the incentives to see that they are minimized (because he captures the benefits from their reduction). Furthermore, these agency costs will be incurred only if the benefits to the owner-manager from their creation are great enough to outweigh them. In our current example these benefits arise from the availability of profitable investments requiring capital investment in excess of the original owner's personal wealth.

In conclusion, finding that agency costs are non-zero (i.e., that there are costs associated with the separation of ownership and control in the corporation) and concluding therefrom that the agency relationship is non-optimal, wasteful or inefficient is equivalent in every sense to comparing a world in which iron ore is a scarce commodity (and therefore costly) to a world in which it is freely available at zero resource cost, and concluding that the first world is "non-optimal" – a perfect example of the fallacy criticized by Coase (1964) and what Demsetz (1969) characterizes as the "Nirvana" form of analysis.[28]

### 2.6. Factors affecting the size of the divergence from ideal maximization

The magnitude of the agency costs discussed above will vary from firm to firm. It will depend on the tastes of managers, the ease with which they can exercise their own preferences as opposed to value maximization in decision making, and the costs of monitoring and bonding activities.[29] The agency costs will also depend upon the cost of measuring the manager's (agent's) performance and evaluating it, the cost of devising and applying an index for compensating the manager which correlates with the owner's (principal's) welfare, and the cost of devising and enforcing specific behavioral rules or policies. Where the manager has less than a controlling interest in the firm, it will also depend upon the market for managers. Competition from other potential managers limits the costs of obtaining managerial services (including the extent to which a given manager can diverge from the idealized solution which would obtain if all monitoring and bonding costs were zero). The size of the divergence (the agency costs) will be directly related to the cost of replacing the manager. If his responsibilities require

---

[28]If we could establish the existence of a feasible set of alternative institutional arrangements which would yield net benefits from the reduction of these costs we could legitimately conclude the agency relationship engendered by the corporation was not Pareto optimal. However, we would then be left with the problem of explaining why these alternative institutional arrangements have not replaced the corporate form of organization.

[29]The monitoring and bonding costs will differ from firm to firm depending on such things as the inherent complexity and geographical dispersion of operations, the attractiveness of perquisites available in the firm (consider the mint), etc.

very little knowledge specialized to the firm, if it is easy to evaluate his performance, and if replacement search costs are modest, the divergence from the ideal will be relatively small and vice versa.

The divergence will also be constrained by the market for the firm itself, i.e., by capital markets. Owners always have the option of selling their firm, either as a unit or piecemeal. Owners of manager-operated firms can and do sample the capital market from time to time. If they discover that the value of the future earnings stream to others is higher than the value of the firm to them given that it is to be manager-operated, they can exercise their right to sell. It is conceivable that other owners could be more efficient at monitoring or even that a single individual with appropriate managerial talents and with sufficiently large personal wealth would elect to buy the firm. In this latter case the purchase by such a single individual would completely eliminate the agency costs. If there were a number of such potential owner–manager purchasers (all with talents and tastes identical to the current manager) the owners would receive in the sale price of the firm the full value of the residual claimant rights including the capital value of the eliminated agency costs plus the value of the managerial rights.

*Monopoly, competition and managerial behavior.* It is frequently argued that the existence of competition in product (and factor) markets will constrain the behavior of managers to idealized value maximization, i.e., that monopoly in product (or monopsony in factor) markets will permit larger divergences from value maximization.[30] Our analysis does not support this hypothesis. The owners of a firm with monopoly power have the same incentives to limit divergences of the manager from value maximization (i.e., the ability to increase their wealth) as do the owners of competitive firms. Furthermore, competition in the market for managers will generally make it unnecessary for the owners to share rents with the manager. The owners of a monopoly firm need only pay the supply price for a manager.

Since the owner of a monopoly has the same wealth incentives to minimize managerial costs as would the owner of a competitive firm, both will undertake that level of monitoring which equates the marginal cost of monitoring to the

---

[30]'Where competitors are numerous and entry is easy, persistent departures from profit maximizing behavior inexorably leads to extinction. Economic natural selection holds the stage. In these circumstances, the behavior of the individual units that constitute the supply side of the product market is essentially routine and uninteresting and economists can confidently predict industry behavior without being explicitly concerned with the behavior of these individual units.

When the conditions of competition are relaxed, however, the opportunity set of the firm is expanded. In this case, the behavior of the firm as a distinct operating unit is of separate interest. Both for purposes of interpreting particular behavior within the firm as well as for predicting responses of the industry aggregate, it may be necessary to identify the factors that influence the firm's choices within this expanded opportunity set and embed these in a formal model.' [Williamson (1964, p. 2)]

marginal wealth increment from reduced consumption of perquisites by the manager. Thus, the existence of monopoly will not increase agency costs.

Furthermore the existence of competition in product and factor markets will not eliminate the agency costs due to managerial control problems as has often been asserted [cf. Friedman (1970)]. If my competitors all incur agency costs equal to or greater than mine I will not be eliminated from the market by their competition.

The existence and size of the agency costs depends on the nature of the monitoring costs, the tastes of managers for non-pecuniary benefits and the supply of potential managers who are capable of financing the entire venture out of their personal wealth. If monitoring costs are zero, agency costs will be zero or if there are enough 100 percent owner–managers available to own and run all the firms in an industry (competitive or not) then agency costs in that industry will also be zero.[31]

## 3. Some unanswered questions regarding the existence of the corporate form

### 3.1. The question

The analysis to this point has left us with a basic puzzle: Why, given the existence of positive costs of the agency relationship, do we find the usual corporate form of organization with widely diffuse ownership so widely prevalent? If one takes seriously much of the literature regarding the "discretionary" power held by managers of large corporations, it is difficult to understand the historical fact of enormous growth in equity in such organizations, not only in the United States, but throughout the world. Paraphrasing Alchian (1968): How does it happen that millions of individuals are willing to turn over a significant fraction of their wealth to organizations run by managers who have so little interest in their welfare? What is even more remarkable, why are they willing to make these commitments purely as residual claimants, i.e., on the anticipation that managers will operate the firm so that there will be earnings which accrue to the stockholders?

There is certainly no lack of alternative ways that individuals might invest, including entirely different forms of organizations. Even if consideration is limited to corporate organizations, there are clearly alternative ways capital might be raised, i.e., through fixed claims of various sorts, bonds, notes, mortgages, etc. Moreover, the corporate income tax seems to favor the use of fixed claims since interest is treated as a tax deductible expense. Those who assert that managers do not behave in the interest or stockholders have generally not addressed a very important question: Why, if non-manager-owned shares have

---

[31]Assuming there are no special tax benefits to ownership nor utility of ownership other than that derived from the direct wealth effects of ownership such as might be true for professional sports teams, race horse stables, firms which carry the family name, etc.

such a serious deficiency, have they not long since been driven out by fixed claims?[32]

### 3.2. Some alternative explanations of the ownership structure of the firm

*The role of limited liability.* Manne (1967) and Alchian and Demsetz (1972) argue that one of the attractive features of the corporate form vis-a-vis individual proprietorships or partnerships is the limited liability feature of equity claims in corporations. Without this provision each and every investor purchasing one or more shares of a corporation would be potentially liable to the full extent of his personal wealth for the debts of the corporation. Few individuals would find this a desirable risk to accept and the major benefits to be obtained from risk reduction through diversification would be to a large extent unobtainable. This argument, however, is incomplete since limited liability does not eliminate the basic risk, it merely shifts it. The argument must rest ultimately on transactions costs. If all stockholders of GM were liable for GM's debts, the maximum liability for an individual shareholder would be greater than it would be if his shares had limited liability. However, given that many other stockholder's also existed and that each was liable for the unpaid claims in proportion to his ownership it is highly unlikely that the maximum payment each would have to make would be large in the event of GM's bankruptcy since the total wealth of those stockholders would also be large. However, the existence of unlimited liability would impose incentives for each shareholder to keep track of both the liabilities of GM and the wealth of the other GM owners. It is easily conceivable that the costs of so doing would, in the aggregate, be much higher than simply paying a premium in the form of higher interest rates to the creditors of GM in return for their acceptance of a contract which grants limited liability to the shareholders. The creditors would then bear the risk of any non-payment of debts in the event of GM's bankruptcy.

It is also not generally recognized that limited liability is merely a necessary condition for explaining the magnitude of the reliance on equities, not a sufficient condition. Ordinary debt also carries limited liability.[33] If limited liability is all that is required, why don't we observe large corporations, individually owned, with a tiny fraction of the capital supplied by the entrepreneur,

---

[32]Marris (1964, pp. 7–9) is the exception, although he argues that there exists some 'maximum leverage point' beyond which the chances of 'insolvency' are in some undefined sense too high.

[33]By limited liability we mean the same conditions that apply to common stock. Subordinated debt or preferred stock could be constructed which carried with it liability provisions; i.e., if the corporation's assets were insufficient at some point to pay off all prior claims (such as trade credit, accrued wages, senior debt, etc.) and if the personal resources of the 'equity' holders were also insufficient to cover these claims the holders of this 'debt' would be subject to assessments beyond the face value of their claim (assessments which might be limited or unlimited in amount).

and the rest simply borrowed[34] At first this question seems silly to many people (as does the question regarding why firms would ever issue debt or preferred stock under conditions where there are no tax benefits obtained from the treatment of interest or preferred dividend payments[35]). We have found that oftentimes this question is misinterpreted to be one regarding why firms obtain capital. The issue is not why they obtain capital, but why they obtain it through the particular forms we have observed for such long periods of time. The fact is that no well articulated answer to this question currently exists in the literature of either finance or economics.

*The "irrelevance" of capital structure.* In their pathbreaking article on the cost of capital, Modigliani and Miller (1958) demonstrated that in the absence of bankruptcy costs and tax subsidies on the payment of interest the value of the firm is independent of the financial structure. They later (1963) demonstrated that the existence of tax subsidies on interest payments would cause the value of the firm to rise with the amount of debt financing by the amount of the capitalized value of the tax subsidy. But this line of argument implies that the firm should be financed almost entirely with debt. Realizing the inconsistence with observed behavior Modigliani and Miller (1963, p. 442) comment:

> "it may be useful to remind readers once again that the existence of a tax advantage for debt financing ... does not necessarily mean that corporations should at all times seek to use the maximum amount of debt in their capital structures. ... there are as we pointed out, limitations imposed by lenders ... as well as many other dimensions (and kinds of costs) in real-world problems of financial strategy which are not fully comprehended within the framework of static equilibrium models, either our own or those of the traditional variety. These additional considerations, which are typically grouped under the rubric of 'the need for preserving flexibility',

[34]Alchian-Demsetz (1972, p. 709) argue that one can explain the existence of both bonds and stock in the ownership structure of firms as the result of differing expectations regarding the outcomes to the firm. They argue that bonds are created and sold to 'pessimists' and stocks with a residual claim with no upper bound are sold to 'optimists'.

As long as capital markets are perfect with no taxes or transactions costs and individual investors can issue claims on distributions of outcomes on the same terms as firms, such actions on the part of firms cannot affect their values. The reason is simple. Suppose such 'pessimists' did exist and yet the firm issues only equity claims. The demand for those equity claims would reflect the fact that the individual purchaser could on his own account issue 'bonds' with a limited and prior claim on the distribution of outcomes on the equity which is exactly the same as that which the firm could issue. Similarly, investors could easily unlever any position by simply buying a proportional claim on both the bonds and stocks of a levered firm. Therefore, a levered firm could not sell at a different price than an unlevered firm solely because of the existence of such differential expectations. See Fama and Miller (1972, ch. 4) for an excellent exposition of these issues.

[35]Corporations did use both prior to the institution of the corporate income tax in the U.S. and preferred dividends have, with minor exceptions, never been tax deductible.

will normally imply the maintenance by the corporation of a substantial reserve of untapped borrowing power."

Modigliani and Miller are essentially left without a theory of the determination of the optimal capital structure, and Fama and Miller (1972, p. 173) commenting on the same issue reiterate this conclusion:

"And we must admit that at this point there is little in the way of convincing research, either theoretical or empirical, that explains the amounts of debt that firms do decide to have in their capital structure."

The Modigliani–Miller theorem is based on the assumption that the probability distribution of the cash flows to the firm is independent of the capital structure. It is now recognized that the existence of positive costs associated with bankruptcy and the presence of tax subsidies on corporate interest payments will invalidate this irrelevance theorem precisely because the probability distribution of future cash flows changes as the probability of the incurrence of the bankruptcy costs changes, i.e., as the ratio of debt to equity rises. We believe the existence of agency costs provide stronger reasons for arguing that the probability distribution of future cash flows is *not* independent of the capital or ownership structure.

While the introduction of bankruptcy costs in the presence of tax subsidies leads to a theory which defines an optimal capital structure,[36] we argue that this theory is seriously incomplete since it implies that no debt should ever be used in the absence of tax subsidies if bankruptcy costs are positive. Since we know debt was commonly used prior to the existence of the current tax subsidies on interest payments this theory does not capture what must be some important determinants of the corporate capital structure.

In addition, neither bankruptcy costs nor the existence of tax subsidies can explain the use of preferred stock or warrnts which have no tax advantages, and there is no theory which tells us anything about what determines the fraction of equity claims held by insiders as opposed to outsiders which our analysis in section 2 indicates is so important. We return to these issues later after analyzing in detail the factors affecting the agency costs associated with debt.

## 4. The agency costs of debt

In general if the agency costs engendered by the existence of outside owners are positive it will pay the absentee owner (i.e., shareholders) to sell out to an owner–manager who can avoid these costs.[37] This could be accomplished in principle by having the manager become the sole equity holder by repurchasing

---

[36]See Kraus and Litzenberger (1972) and Lloyd-Davies (1975).

[37]And if there is competitive bidding for the firm from potential owner-managers the absentee owner will capture the capitalized value of these agency costs.

all of the outside equity claims with funds obtained through the issuance of limited liability debt claims and the use of his own personal wealth. This single-owner corporation would not suffer the agency costs associated with outside equity. Therefore there must be some compelling reasons why we find the diffuse-owner corporate firm financed by equity claims so prevalent as an organizational form.

An ingenious entrepreneur eager to expand, has open to him the opportunity to design a whole hierarchy of fixed claims on assets and earnings, with premiums paid for different levels of risk.[38] Why don't we observe large corporations individually owned with a tiny fraction of the capital supplied by the entrepreneur in return for 100 percent of the equity and the rest simply borrowed? We believe there are a number of reasons: (1) the incentive effects associated with highly leveraged firms, (2) the monitoring costs these incentive effects engender, and (3) bankruptcy costs. Furthermore, all of these costs are simply particular aspects of the agency costs associated with the existence of debt claims on the firm.

### 4.1. The incentive effects associated with debt

We don't find many large firms financed almost entirely with debt type claims (i.e., non-residual claims) because of the effect such a financial structure would have on the owner–manager's behavior. Potential creditors will not loan $100,000,000 to a firm in which the entrepreneur has an investment of $10,000. With that financial structure the owner–manager will have a strong incentive to engage in activities (investments) which promise very high payoffs if successful even if they have a very low probability of success. If they turn out well, he captures most of the gains, if they turn out badly, the creditors bear most of the costs.[39]

To illustrate the incentive effects associated with the existence of debt and to provide a framework within which we can discuss the effects of monitoring and bonding costs, wealth transfers, and the incidence of agency costs, we again consider a simple situation. Assume we have a manager-owned firm with no debt

[38]The spectrum of claims which firms can issue is far more diverse than is suggested by our two-way classification – fixed vs. residual. There are convertible bonds, equipment trust certificates, debentures, revenue bonds, warrants, etc. Different bond issues can contain different subordination provisions with respect to assets and interest. They can be callable or non-callable. Preferred stocks can be 'preferred' in a variety of dimensions and contain a variety of subordination stipulations. In the abstract, we can imagine firms issuing claims contingent on a literally infinite variety of states of the world such as those considered in the literature on the time–state-preference models of Arrow (1964), Debreu (1959) and Hirshleifer (1970).

[39]An apt analogy is the way one would play poker on money borrowed at a fixed interest rate, with one's own liability limited to some very small stake. Fama and Miller (1972, pp. 179–180) also discuss and provide a numerical example of an investment decision which illustrates very nicley the potential inconsistency between the interests of bondholders and stockholders.

outstanding in a world in which there are no taxes. The firm has the opportunity to take one of two mutually exclusive equal cost investment opportunities, each of which yields a random payoff, $\tilde{X}_j$, $T$ periods in the future ($j = 1, 2$). Production and monitoring activities take place continuously between time 0 and time $T$, and markets in which the claims on the firm can be traded are open continuously over this period. After time $T$ the firm has no productive activities so the payoff $\tilde{X}_j$ includes the distribution of all remaining assets. For simplicity, we assume that the two distributions are log-normally distributed and have the same expected total payoff, $E(\tilde{X})$, where $\tilde{X}$ is defined as the logarithm of the final payoff. The distributions differ only by their variances with $\sigma_1^2 < \sigma_2^2$. The systematic or covariance risk of each of the distributions, $\beta_j$, in the Sharpe (1964) – Lintner (1965) capital asset pricing model, is assumed to be identical. Assuming that asset prices are determined according to the capital asset pricing model, the preceding assumptions imply that the total market value of each of these distributions is identical, and we represent this value by $V$.

If the owner–manager has the right to decide which investment program to take, and if after he decides this he has the opportunity to sell part or all of his claims on the outcomes in the form of either debt or equity, he will be indifferent between the two investments.[40]

However, if the owner has the opportunity to *first* issue debt, then to decide which of the investments to take, and then to sell all or part of his remaining equity claim on the market, he will not be indifferent between the two investments. The reason is that by promising to take the low variance project, selling bonds and then taking the high variance project he can transfer wealth from the (naive) bondholders to himself as equity holder.

Let $X^*$ be the amount of the "fixed" claim in the form of a non-coupon bearing bond sold to the bondholders such that the total payoff to them, $R_j$ ($j = 1, 2$, denotes the distribution the manager chooses), is

$$
\begin{aligned}
R_j &= X^*, \quad \text{if} \quad \tilde{X}_j \geqq X^*, \\
    &= X_j, \quad \text{if} \quad \tilde{X}_j \leq X^*.
\end{aligned}
$$

Let $B_1$ be the current market value of bondholder claims if investment 1 is taken, and let $B_2$ be the current market value of bondholders claims if investment 2 is taken. Since in this example the total value of the firm, $V$, is independent of the investment choice and also of the financing decision we can use the Black–Scholes (1973) option pricing model to determine the values of the debt, $B_j$, and equity, $S_j$, under each of the choices.[41]

---

[40]The portfolio diversification issues facing the owner-manager are brought into the analysis in section 5 below.

[41]See Smith (1976) for a review of this option pricing literature and its applications and Galai and Masulis (1976) who apply the option pricing model to mergers, and corporate investment decisions.

Black–Scholes derive the solution for the value of a European call option (one which can be exercised only at the maturity date) and argue that the resulting option pricing equation can be used to determine the value of the equity claim on a levered firm. That is the stockholders in such a firm can be viewed as holding a European call option on the total value of the firm with exercise price equal to $X^*$ (the face value of the debt), exercisable at the maturity date of the debt issue. More simply, the stockholders have the right to buy the firm back from the bondholders for a price of $X^*$ at time $T$. Merton (1973, 1974) shows that as the variance of the outcome distribution rises the value of the stock (i.e., call option) rises, and since our two distributions differ only in their variances, $\sigma_2^2 < \sigma_1^2$, the equity value $S_1$ is less than $S_2$. This implies $B_1 > B_2$, since $B_1 = V - S_1$ and $B_2 = V - S_2$.

Now if the owner–manager could sell bonds with face value $X^*$ under the conditions that the potential bondholders believed this to be a claim on distribution 1, he would receive a price of $B_1$. After selling the bonds, his equity interest in distribution 1 would have value $S_1$. But we know $S_2$ is greater than $S_1$ and thus the manager can make himself better off by changing the investment to take the higher variance distribution 2, thereby redistributing wealth from the bondholders to himself. All this assumes of course that the bondholders could not prevent him from changing the investment program. *If the bondholders cannot do so, and if they perceive that the manager has the opportunity to take distribution 2 they will pay the manager only $B_2$ for the claim $X^*$, realizing that his maximizing behavior will lead him to choose distribution 2.* In this event there is no redistribution of wealth between bondholders and stockholders (and in general with rational expectations there never will be) and no welfare loss. It is easy to construct a case, however, in which these incentive effects do generate real costs.

Let cash flow distribution 2 in the previous example have an expected value, $E(X_2)$, which is lower than that of distribution 1. Then we know that $V_1 > V_2$, and if $\Delta V$, which is given by

$$\Delta V = V_1 - V_2 = (S_1 - S_2) + (B_1 - B_2),$$

is sufficiently small relative to the reduction in the value of the bonds the value of the stock will increase.[42] Rearranging the expression for $\Delta V$ we see that the

---

[42]While we used the option pricing model above to motivate the discussion and provide some intuitive understanding of the incentives facing the equity holders, the option pricing solutions of Black and Scholes (1973) do not apply when incentive effects cause $V$ to be a function of the debt/equity ratio as it is in general and in this example. Long (1974) points out this difficulty with respect to the usefulness of the model in the context of tax subsidies on interest and bankruptcy cost. The results of Merton (1974) and Galai and Masulis (1976) must be interpreted with care since the solutions are strictly incorrect in the context of tax subsidies and/or agency costs.

difference between the equity values for the two investments is given by

$$S_2 - S_1 = (B_1 - B_2) - (V_1 - V_2),$$

and the first term on the RHS, $B_1 - B_2$, is the amount of wealth "transferred" from the bondholders and $V_1 - V_2$ is the reduction in overall firm value. Since we know $B_1 > B_2$, $S_2 - S_1$ can be positive even though the reduction in the value of the firm, $V_1 - V_2$, is positive.[43] Again, the bondholders will not actually lose as long as they accurately perceive the motivation of the equity owning manager and his opportunity to take project 2. They will presume he will take investment 2, and hence will pay no more than $B_2$ for the bonds when they are issued.

In this simple example the reduced value of the firm, $V_1 - V_2$, is the agency cost engendered by the issuance of debt[44] and it is borne by the owner–manager. If he could finance the project out of his personal wealth, he would clearly choose project 1 since its investment outlay was assumed equal to that of project 2 and its market value, $V_1$, was greater. This wealth loss, $V_1 - V_2$, is the "residual loss" portion of what we have defined as agency costs and it is generated by the cooperation required to raise the funds to make the investment. Another important part of the agency costs are monitoring and bonding costs and we now consider their role.

### 4.2. The role of monitoring and bonding costs

In principle it would be possible for the bondholders, by the inclusion of various covenants in the indenture provisions, to limit the managerial behavior

---

[43]The numerical example of Fama and Miller (1972, pp. 179-180) is a close representation of this case in a two-period state model. However, they go on to make the following statement on p. 180:

> 'From a practical viewpoint, however, situations of potential conflict between bondholders and shareholders in the application of the market value rule are probably unimportant. In general, investment opportunities that increase a firm's market value by more than their cost both increase the value of the firm's shares and strengthen the firm's future ability to meet its current bond commitments.'

This first issue regarding the importance of the conflict of interest between bondholders and stockholders is an empirical one, and the last statement is incomplete – in some circumstances the equity holders could benefit from projects whose net effect was to reduce the total value of the firm as they and we have illustrated. The issue cannot be brushed aside so easily.

[44]Myers (1975) points out another serious incentive effect on managerial decisions of the existence of debt which does not occur in our simple single decision world. He shows that if the firm has the option to take future investment opportunities the existence of debt which matures after the options must be taken will cause the firm (using an equity value maximizing investment rule) to refuse to take some otherwise profitable projects because they would benefit only the bondholders and not the equity holders. This will (in the absence of tax subsidies to debt) cause the value of the firm to fall. Thus (although he doesn't use the term) these incentive effects also contribute to the agency costs of debt in a manner perfectly consistent with the examples discussed in the text.

which results in reductions in the value of the bonds. Provisions which impose constraints on management's decisions regarding such things as dividends, future debt issues,[45] and maintenance of working capital are not uncommon in bond issues.[46] To completely protect the bondholders from the incentive effects, these provisions would have to be incredibly detailed and cover most operating aspects of the enterprise including limitations on the riskiness of the projects undertaken. The costs involved in writing such provisions, the costs of enforcing them and the reduced profitability of the firm (induced because the covenants occasionally limit management's ability to take optimal actions on certain issues) would likely be non-trivial. In fact, since management is a continuous decision making process it will be almost impossible to completely specify such conditions without having the bondholders actually perform the management function. All costs associated with such covenants are what we mean by monitoring costs.

The bondholders will have incentives to engage in the writing of such covenants and in monitoring the actions of the manager to the point where the "nominal" marginal cost to them of such activities is just equal to the marginal benefits they perceive from engaging in them. We use the word nominal here because debtholders will not in fact bear these costs. As long as they recognize their existence, they will take them into account in deciding the price they will pay for any given debt claim,[47] and therefore the seller of the claim (the owner) will bear the costs just as in the equity case discussed in section 2.

In addition the manager has incentives to take into account the costs imposed on the firm by covenants in the debt agreement which directly affect the future cash flows of the firm since they reduce the market value of his claims. Because both the external and internal monitoring costs are imposed on the owner-manager it is in his interest to see that the monitoring is performed in the lowest cost way. Suppose, for example, that the bondholders (or outside equity holders) would find it worthwhile to produce detailed financial statements such as those contained in the usual published accounting reports as a means of monitoring the manager. If the manager himself can produce such information at lower costs than they (perhaps because he is already collecting much of the data they desire for his own internal decision making purposes), it would pay him to agree in advance to incur the cost of providing such reports and to have their

[45]Black–Scholes (1973) discuss ways in which dividend and future financing policy can redistribute wealth between classes of claimants on the firm.

[46]Black, Miller and Posner (1974) discuss many of these issues with particular reference to the government regulation of bank holding companies.

[47]In other words, these costs will be taken into account in determing the yield to maturity on the issue. For an examination of the effects of such enforcement costs on the nominal interest rates in the consumer small loan market, see Benston (1977).

accuracy testified to by an independent outside auditor. This is an example of what we refer to as bonding costs.[48,49]

### 4.3. Bankruptcy and reorganization costs

We argue in section 5 that as the debt in the capital structure increases beyond some point the marginal agency costs of debt begin to dominate the marginal

[48]To illustrate the fact that it will sometimes pay the manager to incur 'bonding' costs to guarantee the bondholders that he will not deviate from his promised behavior let us suppose that for an expenditure of $b$ of the firm's resources he can guarantee that project 1 will be chosen. If he spends these resources and takes project 1 the value of the firm will be $V_1 - b$ and clearly as long as $(V_1 - b) > V_2$, or alternatively $(V_1 - V_2) > b$ he will be better off, since his wealth will be equal to the value of the firm minus the required investment, $I$ (which we assumed for simplicity to be identical for the two projects).

On the other hand, to prove that the owner–manager prefers the lowest cost solution to the conflict let us assume he can write a covenant into the bond issue which will allow the bondholders to prevent him from taking project 2, if they incur monitoring costs of $m$, where $m < b$. If he does this his wealth will be higher by the amount $b - m$. To see this note that if the bond market is competitive and makes unbiased estimates, potential bondholders will be indifferent between:

(i) a claim $X^*$ with no covenant (and no guarantees from management) at a price of $B_2$,
(ii) a claim $X^*$ with no covenant (and guarantees from management, through bonding expenditures by the firm of $b$, that project 1 will be taken) at a price of $B_1$, and
(iii) a claim $X^*$ with a covenant and the opportunity to spend $m$ on monitoring (to guarantee project 1 will be taken) at a price of $B_1 - m$.

The bondholders will realize that (i) represents in fact a claim on project 2 and that (ii) and (iii) represent a claim on project 1 and are thus indifferent between the three options at the specified prices. The owner–manager, however, will not be indifferent between incurring the bonding costs, $b$, directly, or including the covenant in the bond indenture and letting the bondholders spend $m$ to guarantee that he take project 1. His wealth in the two cases will be given by the value of his equity plus the proceeds of the bond issue less the required investment, and if $m < b < V_1 - V_2$, then his post-investment-financing wealth, $W$, for the three options will be such that $W_I < W_{II} < W_{III}$. Therefore, since it would increase his wealth, he would voluntarily include the covenant in the bond issue and let the bondholders monitor.

[49]We mention, without going into the problem in detail, that similar to the case in which the outside equity holders are allowed to monitor the manager-owner, the agency relationship between the bondholders and stockholders has a symmetry if the rights of the bondholders to limit actions of the manager are not perfectly spelled out. Suppose the bondholders, by spending sufficiently large amounts of resources, could force management to take actions which would transfer wealth from the equity holder to the bondholders (by taking sufficiently less risky projects). One can easily construct situations where such actions could make the bondholders better off, hurt the equity holders and actually lower the total value of the firm. Given the nature of the debt contract the original owner-manager might maximize his wealth in such a situation by selling off the equity and keeping the bonds as his 'owner's' interest. If the nature of the bond contract is given, this may well be an inefficient solution since the total agency costs (i.e., the sum of monitoring and value loss) could easily be higher than the alternative solution. However, if the owner-manager could strictly limit the rights of the bondholders (perhaps by inclusion of a provision which expressly reserves all rights not specifically granted to the bondholder for the equity holder), he would find it in his interest to establish the efficient contractual arrangement since by minimizing the agency costs he would be maximizing his wealth. These issues involve the fundamental nature of contracts and for now we simply assume that the 'bondholders' rights are strictly limited and unambiguous and all rights not specifically granted them are reserved for the 'stockholders'; a situation descriptive of actual institutional arrangements. This allows us to avoid the incentive effects associated with 'bondholders' potentially exploiting 'stockholders'.

agency costs of outside equity and the result of this is the generally observed phenomenon of the simultaneous use of both debt snd outside equity. Before considering these issues, however, we consider here the third major component of the agency costs of debt which helps to explain why debt doesn't completely dominate capital structures – the existence of bankruptcy and reorganization costs.

It is important to emphasize that bankruptcy and liquidation are very different events. The legal definition of bankruptcy is difficult to specify precisely. In general, it occurs when the firm cannot meet a current payment on a debt obligation,[50] or one or more of the other indenture provisions providing for bankruptcy is violated by the firm. In this event the stockholders have lost all claims on the firm,[51] and the remaining loss, the difference between the face value of the fixed claims and the market value of the firm, is borne by the debtholders. Liquidation of the firm's assets will occur only if the market value of the future cash flows generated by the firm is less than the opportunity cost of the assets, i.e., the sum of the values which could be realized if the assets were sold piecemeal.

If there were no costs associated with the event called bankruptcy the total market value of the firm would not be affected by increasing the probability of its incurrence. However, it is costly, if not impossible, to write contracts representing claims on a firm which clearly delineate the rights of holders for all possible contingencies. Thus even if there were no adverse incentive effects in expanding fixed claims relative to equity in a firm, the use of such fixed claims would be constrained by the costs inherent in defining and enforcing those claims. Firms incur obligations daily to suppliers, to employees, to different classes of investors, etc. So long as the firm is prospering, the adjudication of claims is seldom a problem. When the firm has difficulty meeting some of its obligations, however, the issue of the priority of those claims can pose serious problems. This is most obvious in the extreme case where the firm is forced into bankruptcy. If bankruptcy were costless, the reorganization would be accompanied by an adjustment of the claims of various parties and the business, could, if that proved to be in the interest of the claimants, simply go on (although perhaps under new management).[52]

---

[50]If the firm were allowed to sell assets to meet a current debt obligation, bankruptcy would occur when the total market value of the future cash flows expected to be generated by the firm is less than the value of a current payment on a debt obligation. Many bond indentures do not, however, allow for the sale of assets to meet debt obligations.

[51]We have been told that while this is true in principle, the actual behavior of the courts appears to frequently involve the provision of some settlement to the common stockholders even when the assets of the company are not sufficient to cover the claims of the creditors.

[52]If under bankruptcy the bondholders have the right to fire the management, the management will have some incentives to avoid taking actions which increase the probability of this event (even if it is in the best interest of the equity holders) if they (the management) are earning rents or if they have human capital specialized to this firm or if they face large adjustment costs in finding new employment. A detailed examination of this issue involves the value of the control rights (the rights to hire and fire the manager) and we leave it to a subsequent paper.

In practice, bankruptcy is not costless, but generally involves an adjudication process which itself consumes a fraction of the remaining value of the assets of the firm. Thus the cost of bankruptcy will be of concern to potential buyers of fixed claims in the firm since their existence will reduce the payoffs to them in the event of bankruptcy. These are examples of the agency costs of cooperative efforts among individuals (although in this case perhaps "non-cooperative" would be a better term). The price buyers will be willing to pay for fixed claims will thus be inversely related to the probability of the incurrence of these costs i.e., to the probability of bankruptcy. Using a variant of the argument employed above for monitoring costs, it can be shown that the total value of the firm will fall, and the owner–manager equity holder will bear the entire wealth effect of the bankruptcy costs as long as potential bondholders make unbiased estimates of their magnitude at the time they initially purchase bonds.[53]

Empirical studies of the magnitude of bankruptcy costs are almost non-existent. Warner (1975) in a study of 11 railroad bankruptcies between 1930 and 1955 estimates the average costs of bankruptcy[54] as a fraction of the value of the firm three years prior to bankruptcy to be 2.5% (with a range of 0.4% to 5.9%). The average dollar costs were $1.88 million. Both of these measures seem remarkably small and are consistent with our belief that bankruptcy costs themselves are unlikely to be the major determinant of corporate capital structures. It is also interesting to note that the annual amount of defaulted funds has fallen significantly since 1940. [See Atkinson (1967).] One possible explanation for this phenomena is that firms are using mergers to avoid the costs of bankruptcy. This hypothesis seems even more reasonable, if, as is frequently the case, reorganization costs represent only a fraction of the costs associated with bankruptcy.

In general the revenues or the operating costs of the firm are not independent of the probability of bankruptcy and thus the capital structure of the firm. As the probability of bankruptcy increases, both the operating costs and the revenues of the firm are adversely affected, and some of these costs can be avoided by merger. For example, a firm with a high probability of bankruptcy will also find that it must pay higher salaries to induce executives to accept the higher risk of unemployment. Furthermore, in certain kinds of durable goods industries the demand function for the firm's product will not be independent of the probability of bankruptcy. The computer industry is a good example. There, the buyer's welfare is dependent to a significant extent on the ability to maintain the equipment, and on continuous hardware and software development. Furthermore, the owner of a large computer often receives benefits from the software

---

[53]Kraus and Litzenberger (1972) and Lloyd-Davies (1975) demonstrate that the total value of the firm will be reduced by these costs.

[54]These include only payments to all parties for legal fees, professional services, trustees' fees and filing fees. They do not include the costs of management time or changes in cash flows due to shifts in the firm's demand or cost functions discussed below.

developments of other users. Thus if the manufacturer leaves the business or loses his software support and development experts because of financial diffi-culties, the value of the equipment to his users will decline. The buyers of such services have a continuing interest in the manufacturer's viability not unlike that of a bondholder, except that their benefits come in the form of continuing services at lower cost rather than principle and interest payments. Service facilities and spare parts for automobiles and machinery are other examples.

In summary then the agency costs associated with debt[55] consist of:

(1) the opportunity wealth loss caused by the impact of debt on the investment decisions of the firm,

(2) the monitoring and bonding expenditures by the bondholders and the owner–manager (i.e., the firm),

(3) the bankruptcy and reorganization costs.

## 4.4. Why are the agency costs of debt incurred?

We have argued that the owner–manager bears the entire wealth effects of the agency costs of debt and he captures the gains from reducing them. Thus, the agency costs associated with debt discussed above will tend, in the absence of other mitigating factors, to discourage the use of corporate debt. What are the factors that encourage its use?

One factor is the tax subsidy on interest payments. (This will not explain preferred stock where dividends are not tax deductible.)[56] Modigliani and Miller (1963) originally demonstrated that the use of riskless perpetual debt will increase the total value of the firm (ignoring the agency costs) by an amount equal to $\tau B$, where $\tau$ is the marginal and average corporate tax rate and $B$ is the market value of the debt. Fama and Miller (1972, ch. 4) demonstrate that for the case of risky debt the value of the firm will increase by the market value of the (uncertain) tax subsidy on the interest payments. Again, these gains will accrue entirely to

[55]Which, incidentally, exist only when the debt has some probability of default.

[56]Our theory is capable of explaining why in the absence of the tax subsidy on interest pay-ments, we would expect to find firms using both debt and preferred stocks – a problem which has long puzzled at least one of the authors. If preferred stock has all the characteristics of debt except for the fact that its holders cannot put the firm into bankruptcy in the event of nonpayment of the preferred dividends, then the agency costs associated with the issuance of preferred stock will be lower than those associated with debt by the present value of the bank-ruptcy costs.

However, these lower agency costs of preferred stock exist only over some range if as the amount of such stock rises the incentive effects caused by their existence impose value reduc-tions which are larger than that caused by debt (including the bankruptcy costs of debt). There are two reasons for this. First, the equity holder's claims can be eliminated by the debtholders in the event of bankruptcy, and second, the debtholders have the right to fire the management in the event of bankruptcy. Both of these will tend to become more important as an advantage to the issuance of debt as we compare situations with large amounts of preferred stock to equivalent situations with large amounts of debt because they will tend to reduce the incentive effects of large amounts of preferred stock.

the equity and will provide an incentive to utilize debt to the point where the marginal wealth benefits of the tax subsidy are just equal to the marginal wealth effects of the agency costs discussed above.

However, even in the absence of these tax benefits, debt would be utilized if the ability to exploit potentially profitable investment opportunities is limited by the resources of the owner. If the owner of a project cannot raise capital he will suffer an opportunity loss represented by the increment in value offered to him by the additional investment opportunities. Thus even though he will bear the agency costs from selling debt, he will find it desirable to incur them to obtain additional capital as long as the marginal wealth increments from the new investments projects are greater than the marginal agency costs of debt, and these agency costs are in turn less than those caused by the sale of additional equity discussed in section 2. Furthermore, this solution is optimal from the social viewpoint. However, in the absence of tax subsidies on debt these projects must be unique to this firm[57] or they would be taken by other competitive entrepreneurs (perhaps new ones) who possessed the requisite personal wealth to fully finance the projects[58] and therefore able to avoid the existence of debt or outside equity.

## 5. A theory of the corporate ownership structure

In the previous sections we discussed the nature of agency costs associated with outside claims on the firm – both debt and equity. Our purpose here is to integrate these concepts into the beginnings of a theory of the corporate owner-. ship structure. We use the term "ownership structure" rather than "capital structure" to highlight the fact that the crucial variables to be determined are not just the relative amounts of debt and equity but also the fraction of the equity held by the manager. Thus, for a given size firm we want a theory to determine three variables:[59]

---

[57]One other conditions also has to hold to justify the incurrence of the costs associated with the use of debt or outside equity in our firm. If there are other individuals in the economy who have sufficiently large amounts of personal capital to finance the entire firm, our capital constrained owner can realize the full capital value of his current and prospective projects and avoid the agency costs by simply selling the firm (i.e. the, right to take these projects) to one of these individuals. He will then avoid the wealth losses associated with the agency costs caused by the sale of debt or outside equity. If no such individuals exist, it will pay him (and society) to obtain the additional capital in the debt market. This implies, incidentally, that it is somewhat misleading to speak of the owner-manager as the individual who bears the agency costs. One could argue that it is the project which bears the costs since, if it is not sufficiently profitable to cover all the costs (including the agency costs), it will not be taken. We continue to speak of the owner-manager bearing these costs to emphasize the more correct and important point that he has the incentive to reduce them because, if he does, his wealth will be increased.

[58]We continue to ignore for the moment the additional complicating factor involved with the portfolio decisions of the owner, and the implied acceptance of potentially diversifiable risk by such 100% owners in this example.

[59]We continue to ignore such instruments as convertible bonds and warrants.

$S_i$ : inside equity (held by the manager),

$S_o$ : outside equity (held by anyone outside of the firm),

$B$ : debt (held by anyone outside of the firm).

The total market value of the equity is $S = S_i + S_o$, and the total market value of the firm is $V = S + B$. In addition, we also wish to have a theory which determines the optimal size of the firm, i.e., its level of investment.

## 5.1. *Determination of the optimal ratio of outside equity to debt*

Consider first the determination of the optimal ratio of outside equity to debt, $S_o/B$. To do this let us hold the size of the firm constant. $V$, the actual value of the firm for a given size, will depend on the agency costs incurred, hence we use as our index of size $V^*$, the value of the firm at a given scale when agency costs are zero. For the moment we also hold the amount of outside financing $(B + S_o)$, constant. Given that a specified amount of financing $(B + S_o)$ is to be obtained externally our problem is to determine the optimal fraction $E^* \equiv S_o^*/(B + S_o)$ to be financed with equity.



Fig. 5. Total agency costs, $A_T(E)$, as a function of the ratio of outside equity, to total outside financing, $E \equiv S_o/(B + S_o)$, for a given firm size $V^*$ and given total amounts of outside financing $(B + S_o)$. $A_{S_o}(E) \equiv$ agency costs associated with outside equity, $A_B(E) \equiv$ agency costs associated with debt, $B$. $A_T(E^*) =$ minimum total agency costs at optimal fraction of outside financing $E^*$.

We argued above that: (1) as long as capital markets are efficient (i.e., characterized by rational expectations) the prices of assets such as debt and outside equity will reflect unbiased estimates of the monitoring costs and redistributions which the agency relationship will engender, and (2) the selling owner–manager will bear these agency costs. Thus from the owner–manager's standpoint the optimal proportion of outside funds to be obtained from equity (versus debt) *for a given level of internal equity* is that $E$ which results in minimum total agency costs.

Fig. 5 presents a breakdown of the agency costs into two separate components: Define $A_{S_o}(E)$ as the total agency costs (a function of $E$) associated with the 'exploitation' of the outside equity holders by the owner–manager, and $A_B(E)$ as the total agency costs associated with the presence of debt in the ownership structure. $A_T(E) = A_{S_o}(E) + A_B(E)$ is the total agency cost.

Consider the function $A_{S_o}(E)$. When $E \equiv S_o/(B + S_o)$ is zero, i.e., when there is no outside equity, the manager's incentives to exploit the outside equity is at a minimum (zero) since the changes in the value of the *total* equity are equal to the changes in *his* equity.[60] As $E$ increases to 100 percent his incentives to exploit the outside equity holders increase and hence the agency costs $A_{S_o}(E)$ increase.

The agency costs associated with the existence of debt, $A_B(E)$ are composed mainly of the value reductions in the firm and monitoring costs caused by the manager's incentive to reallocate wealth from the bondholders to himself by increasing the value of his equity claim. They are at a maximum where all outside funds are obtained from debt, i.e., where $S_o = E = 0$. As the amount of debt declines to zero these costs also go to zero because as $E$ goes to 1, his incentive to reallocate wealth from the bondholders to himself falls. These incentives fall for two reasons: (1) the total amount of debt falls, and therefore it is more difficult to reallocate any given amount away from the debtholders, and (2) his share of any reallocation which is accomplished is falling since $S_o$ is rising and therefore $S_i/(S_o + S_i)$, his share of the total equity, is falling.

The curve $A_T(E)$ represents the sum of the agency costs from various combinations of outside equity and debt financing, and as long as $A_{S_o}(E)$ and $A_B(E)$ are

---

[60]Note, however, that even when outsiders own none of the equity the stockholder-manager still has some incentives to engage in activities which yield him non-pecuniary benefits but reduce the value of the firm by more than he personally values the benefits if there is any risky debt outstanding. Any such actions he takes which reduce the value of the firm, $V$, tend to reduce the value of the bonds as well as the value of the equity. Although the option pricing model does not in general apply exactly to the problem of valuing the debt and equity of the firm, it can be useful in obtaining some qualitative insights into matters such as this. In the option pricing model $\partial S/\partial V$ indicates the rate at which the stock value changes per dollar change in the value of the firm (and similarly for $\partial B/\partial V$). Both of these terms are less than unity [cf. Black and Scholes (1973)]. Therefore, any action of the manager which reduces the value of the firm, $V$, tends to reduce the value of both the stock and the bonds, and the larger is the total debt/equity ratio the smaller is the impact of any given change in $V$ on the value of the equity, and therefore, the lower is the cost to him of consuming non-pecuniary benefits.

as we have drawn them the minimum total agency cost for given size firm and outside financing will occur at some point such as $A_T(E^*)$ with a mixture of both debt and equity.[61]

*A caveat.*   Before proceeding further we point out that the issue regarding the exact shapes of the functions drawn in fig. 5 and several others discussed below is essentially an open question at this time. In the end the shape of these functions is a question of fact and can only be settled by empirical evidence. We outline some a priori arguments which we believe lead to some plausible hypotheses about the behavior of the system, but confess that we are far from understanding the many conceptual subtleties of the problem. We are fairly confident of our arguments regarding the signs of the first derivatives of the functions, but the second derivatives are also important to the final solution and much more work (both theoretical and empirical) is required before we can have much confidence regarding these parameters. We anticipate the work of others as well as our own to cast more light on these issues. Moreover, we suspect the results of such efforts will generate revisions to the details of what follows. We believe it is worthwhile to delineate the overall framework in order to demonstrate, if only in a simplified fashion, how the major pieces of the puzzle fit together into a cohesive structure.

### 5.2. Effects of the scale of outside financing

In order to investigate the effects of increasing the amount of outside financing, $B+S_o$, and therefore reducing the amount of equity held by the manager, $S_i$, we continue to hold the scale of the firm, $V^*$, constant. Fig. 6 presents a plot of the agency cost functions, $A_{S_o}(E)$, $A_B(E)$ and $A_T(E) = A_{S_o}(E) + A_B(E)$, for two different levels of outside financing. Define an index of the amount of outside financing to be

$$K = (B+S_o)/V^*,$$

and consider two different possible levels of outside financing $K_o$ and $K_1$ for a given scale of the firm such that $K_o < K_1$.

As the amount of outside equity increases, the owner's fractional claim on the firm, $\alpha$, falls. He will be induced thereby to take additional non-pecuniary benefits out of the firm because his share of the cost falls. This also increases the marginal benefits from monitoring activities and therefore will tend to increase the optimal level of monitoring. Both of these factors will cause the locus of agency costs $A_{S_o}(E; K)$ to shift upward as the fraction of outside financing, $K$,

---

[61]This occurs, of course, not at the intersection of $A_{S_o}(E)$ and $A_B(E)$, but at the point where the absolute value of the slopes of the functions are equal, i.e., where $A'_{S_o}(E) + A'_B(E) = 0$.

increases. This is depicted in fig. 6 by the two curves representing the agency costs of equity, one for the low level of outside financing, $A_{S_0}(E; K_0)$, the other for the high level of outside financing, $A_{S_0}(E; K_1)$. The locus of the latter lies above the former everywhere except at the origin where both are 0.

The agency cost of debt will similarly rise as the amount of outside financing increases. This means that the locus of $A_B(E; K_1)$ for high outside financing, $K_1$,



Fig. 6. Agency cost functions and optimal outside equity as a fraction of total outside financing, $E^*(K)$, for two different levels of outside financing, $K$, for a given size firm, $V^*$: $K_1 > K_0$.

will lie above the locus of $A_B(E; K_0)$ for low outside financing, $K_0$ because the total amount of resources which can be reallocated from bondholders increases as the total amount of debt increases. However, since these costs are zero when the debt is zero for both $K_0$ and $K_1$ the intercepts of the $A_B(E; K)$ curves coincide at the right axis.

The net effect of the increased use of outside financing given the cost functions assumed in fig. 6 is to: (1) increase the total agency costs from $A_T(E^*; K_0)$ to $A_T(E^*; K_1)$, and (2) to increase the optimal fraction of outside funds obtained from the sale of outside equity. We draw these functions for illustration only and are unwilling to speculate at this time on the exact form of $E^*(K)$ which

65

gives the general effects of increasing outside financing on the relative quantities of debt and equity.

The locus of points, $A_T(E^*; K)$ where agency costs are minimized (not drawn in fig. 6), determines $E^*(K)$, the optimal proportions of equity and debt to be used in obtaining outside funds as the fraction of outside funds, $K$, ranges from 0 to 100 percent. The solid line in fig. 7 is a plot of the minimum total agency costs



Fig. 7. Total agency costs as a function of the fraction of the firm financed by outside claims for two firm sizes, $V_1^* > V_0^*$.

as a function of the amount of outside financing for a firm with scale $V_0^*$. The dotted line shows the total agency costs for a larger firm with scale $V_1^* > V_0^*$. That is, we hypothesize that the larger the firm becomes the larger are the total agency costs because it is likely that the monitoring function is inherently more difficult and expensive in a larger organization.

## 5.3. Risk and the demand for outside financing

The model we have used to explain the existence of minority shareholders and debt in the capital structure of corporations implies that the owner–manager, if he resorts to any outside funding, will have his entire wealth invested in the firm. The reason is that he can thereby avoid the agency costs which additional outside funding impose. This suggests he would not resort to outside funding until he had invested 100 percent of his personal wealth in the firm – an implica-

tion which is not consistent with what we generally observe. Most owner-managers hold personal wealth in a variety of forms, and some have only a relatively small fraction of their wealth invested in the corporation they manage.[62] Diversification on the part of owner-managers can be explained by risk aversion and optimal portfolio selection.

If the returns from assets are not perfectly correlated an individual can reduce the riskiness of the returns on his portfolio by dividing his wealth among many different assets, i.e., by diversifying.[63] Thus a manager who invests all of his wealth in a single firm (his own) will generally bear a welfare loss (if he is risk averse) because he is bearing more risk than necessary. He will, of course, be willing to pay something to avoid this risk, and the costs he must bear to accomplish this diversification will be the agency costs outlined above. He will suffer a wealth loss as he reduces his fractional ownership because prospective shareholders and bondholders will take into account the agency costs. Nevertheless, the manager's desire to avoid risk will contribute to his becoming a minority stockholder.

### 5.4. Determination of the optimal amount of outside financing, K*

Assume for the moment that the owner of a project (i.e., the owner of a prospective firm) has enough wealth to finance the entire project himself. The optimal scale of the corporation is then determined by the condition that, $\Delta V - \Delta I = 0$. In general if the returns to the firm are uncertain the owner-manager can increase his welfare by selling off part of the firm either as debt or equity and reinvesting the proceeds in other assets. If he does this with the optimal combination of debt and equity (as in fig. 6) the total wealth reduction he will incur is given by the agency cost function, $A_t(E^*, K; V^*)$ in fig. 7. The functions $A_T(E^*, K; V^*)$ will be S shaped (as drawn) if total agency costs for a given scale of firm increase at an increasing rate at low levels of outside financing, and at a decreasing rate for high levels of outside financing as monitoring imposes more and more constraints on the manager's actions.

Fig. 8 shows marginal agency costs as a function of $K$, the fraction of the firm financed with outside funds assuming the total agency cost function is as plotted in fig. 7, and assuming the scale of the firm is fixed. The demand by the owner-manager for outside financing is shown by the remaining curve in fig. 8. This curve represents the marginal value of the increased diversification which the manager

---

[62]On the average, however, top managers seem to have substantial holdings in absolute dollars. A recent survey by Wytmar (*Wall Street Journal*, August 13, 1974, p. 1) reported that the median value of 826 chief executive officers' stock holdings in their companies at year end 1973 was $557,000 and $1.3 million at year end 1972.

[63]These diversification effects can be substantial. Evans and Archer (1968) show that on the average for New York Stock Exchange securities approximately 55% of the total risk (as measured by standard deviation of portfolio returns) can be eliminated by following a naive strategy of dividing one's assets equally among 40 randomly selected securities.

can obtain by reducing his ownership claims and optimally constructing a diversified portfolio. It is measured by the amount he would pay to be allowed to reduce his ownership claims by a dollar in order to increase his diversification. If the liquidation of some of his holdings also influences the owner–manager's consumption set, the demand function plotted in fig. 8 also incorporates the marginal value of these effects. The intersection of these two schedules determines



Fig. 8. Determination of the optimal amount of outside financing, $K^*$, for a given scale of firm

the optimal fraction of the firm to be held by outsiders and this in turn determines the total agency costs borne by the owner. This solution is Pareto optimal; there is no way to reduce the agency costs without making someone worse off.

## 5.5. Determination of the optimal scale of the firm

While the details of the solution of the optimal scale of the firm are complicated when we allow for the issuance of debt, equity and monitoring and bonding, the general structure of the solution is analogous to the case where monitoring and bonding are allowed for the outside equity example (see fig. 4).

If it is optimal to issue any debt, the expansion path taking full account of such opportunities must lie above the curve $ZG$ in fig. 4. If this new expansion path lies anywhere to the right of the indifference curve passing through point $G$ debt will be used in the optimal financing package. Furthermore, the optimal scale

of the firm will be determined by the point at which this new expansion path touches the highest indifference curve. In this situation the resulting level of the owner–manager's welfare must therefore be higher.

## 6. Qualifications and extensions of the analysis

### 6.1. Multiperiod aspects of the agency problem

We have assumed throughout our analysis that we are dealing only with a single investment-financing decision by the entrepreneur and have ignored the issues associated with the incentives affecting future financing–investment decisions which might arise after the initial set of contracts are consumated between the entrepreneur–manager, outside stockholders and bondholders. These are important issues which are left for future analysis.[64] Their solution will undoubtedly introduce some changes in the conclusions of the single decision analysis. It seems clear for instance that the expectation of future sales of outside equity and debt will change the costs and benefits facing the manager in making decisions which benefit himself at the (short-run) expense of the current bondholders and stockholders. If he develops a reputation for such dealings he can expect this to unfavourably influence the terms at which he can obtain future capital from outside sources. This will tend to increase the benefits associated with "sainthood" and will tend to reduce the size of the agency costs. Given the finite life of any individual, however, such an effect cannot reduce these costs to zero, because at some point these future costs will begin to weigh more heavily on his successors and therefore the relative benefits to him of acting in his own best interests will rise.[65] Furthermore, it will generally be impossible for him to fully guarantee the outside interests that his successor will continue to follow his policies.

### 6.2. The control problem and outside owner's agency costs

The careful reader will notice that nowhere in the analysis thus far have we taken into account many of the details of the relationship between the part owner–manager and the outside stockholders and bondholders. In particular we have assumed that all outside equity is nonvoting. If such equity does have voting rights then the manager will be concerned about the effects on his long-run welfare of reducing his fractional ownership below the point where he loses

---

[64] The recent work of Myers (1975) which views future investment opportunities as options and investigates the incentive effects of the existence of debt in such a world where a sequence of investment decisions is made is another important step in the investigation of the multi-period aspects of the agency problem and the theory of the firm.

[65] Becker and Stigler (1972) analyze a special case of this problem involving the use of non-vested pension rights to help correct for this end game play in the law enforcement area.

effective control of the corporation. That is, below the point where it becomes possible for the outside equity holders to fire him. A complete analysis of this issue will require a careful specification of the contractual rights involved on both sides, the role of the board of directors, and the coordination (agency) costs borne by the stockholders in implementing policy changes. This latter point involves consideration of the distribution of the outside ownership claims. Simply put, forces exist to determine an equilibrium distribution of outside ownership. If the costs of reducing the dispersion of ownership are lower than the benefits to be obtained from reducing the agency costs, it will pay some individual or group of individuals to buy shares in the market to reduce the dispersion of ownership. We occasionally witness these conflicts for control which involve outright market purchases, tender offers and proxy fights. Further analysis of these issues is left to the future.

### 6.3. A note on the existence of inside debt and some conjectures on the use of convertible financial instruments

We have been asked[66] why debt held by the manager (i.e., "inside debt") plays no role in our analysis. We have as yet been unable to incorporate this dimension formally into our analysis in a satisfactory way. The question is a good one and suggests some potentially important extensions of the analysis. For instance, it suggests an inexpensive way for the owner–manager with both equity and debt outstanding to eliminate a large part (perhaps all) of the agency costs of debt. If he binds himself contractually to hold a fraction of the total debt equal to his fractional ownership of the total equity he would have no incentive whatsoever to reallocate wealth from the debt holders to the stockholders. Consider the case where

$$B_i/S_i = B_o/S_o, \tag{4}$$

where $S_i$ and $S_o$ are as defined earlier, $B_i$ is the dollar value of the inside debt held by the owner–manager, and $B_o$ is the debt held by outsiders. In this case if the manager changes the investment policy of the firm to reallocate wealth between the debt and equity holders, the net effect on the total value of his holdings in the firm will be zero. Therefore, his incentives to perform such reallocations are zero.[67]

Why then don't we observe practices or formal contracts which accomplish

[66] By our colleague David Henderson.

[67] This also suggests that *some* outside debt holders can protect themselves from 'exploitation' by the manager by purchasing a fraction of the total equity equal to their fractional ownership of the debt. All debt holders, of course, cannot do this unless the manager does so also. In addition, such an investment rule restricts the portfolio choices of investors and therefore would impose costs if followed rigidly. Thus the agency costs will not be eliminated this way either.

this elimination or reduction of the agency costs of debt? Maybe we do for smaller privately held firms (we haven't attempted to obtain this data), but for large diffuse owner corporations the practice does not seem to be common. One reason for this we believe is that in some respects the claim that the manager holds on the firm in the form of his wage contract has some of the characteristics of debt.[68] If true, this implies that even with zero holdings of formal debt claims he still has positive holdings of a quasi-debt claim and this may accomplish the satisfaction of condition (4). The problem here is that any formal analysis of this issue requires a much deeper understanding of the relationship between formal debt holdings and the wage contract; i.e., how much debt is it equivalent to?

This line of thought also suggests some other interesting issues. Suppose the implicit debt characteristics of the manager's wage contract result in a situation equivalent to

$$B_i/S_i > B_o/S_o.$$

Then he would have incentives to change the operating characteristics of the firm (i.e., reduce the variance of the outcome distribution) to transfer wealth from the stockholders to the debt holders which is the reverse of the situation we examined in section 4. Furthermore, this seems to capture some of the concern often expressed regarding the fact that managers of large publicly held corporations seem to behave in a risk averse way to the detriment of the equity holders. One solution to this would be to establish incentive compensation systems for the manager or to give him stock options which in effect give him a claim on the upper tail of the outcome distribution. This also seems to be a commonly observed phenomenon.

This analysis also suggests some additional issues regarding the costs and benefits associated with the use of more complicated financial claims such as warrants, convertible bonds and convertible preferred stock which we have not formally analyzed as yet. Warrants, convertible bonds and convertible preferred stock have some of the characteristics of non-voting shares although they can be converted into voting shares under some terms. Alchian–Demsetz (1972) provide an interesting analysis regarding the use of non-voting shares. They argue that some shareholders with strong beliefs in the talents and judgements of the manager will want to be protected against the possibility that some other shareholders will take over and limit the actions of the manager (or fire him). Given that the securities exchanges prohibit the use of non-voting shares by listed firms the use of option type securities might be a substitute for these claims.

In addition warrants represents a claim on the upper tail of the distribution of

---

[68]Consider the situation in which the bondholders have the right in the event of bankruptcy to terminate his employment and therefore to terminate the future returns to any specific human capital or rents he may be receiving.

outcomes, and convertible securities can be thought of as securities with non-detachable warrants. It seems that the incentive effects of warrants would tend to offset to some extent the incentive effects of the existence of risky debt because the owner–manager would be sharing part of the proceeds associated with a shift in the distribution of returns with the warrant holders. Thus, we conjecture that potential bondholders will find it attractive to have warrants attached to the risky debt of firms in which it is relatively easy to shift the distribution of outcomes to expand the upper tail of the distribution to transfer wealth from bondholders. It would also then be attractive to the owner–manager because of the reduction in the agency costs which he would bear. This argument also implies that it would make little difference if the warrants were detachable (and therefore saleable separately from the bonds) since their mere existence would reduce the incentives of the manager (or stockholders) to increase the riskiness of the firm (and therefore increase the probability of bankruptcy). Furthermore, the addition of a conversion privilege to fixed claims such as debt or preferred stock would also tend to reduce the incentive effects of the existence of such fixed claims and therefore lower the agency costs associated with them. The theory predicts that these phenomena should be more frequently observed in cases where the incentive effects of such fixed claims are high than when they are low.

### 6.4. Monitoring and the social product of security analysts

One of the areas in which further analysis is likely to lead to high payoffs is that of monitoring. We currently have little which could be glorified by the title of a "Theory of Monitoring" and yet this is a crucial building block of the analysis. We would expect monitoring activities to become specialized to those institutions and individuals who possess comparative advantages in these activities. One of the groups who seem to play a large role in these activities is composed of the security analysts employed by institutional investors, brokers and investment advisory services as well as the analysis performed by individual investors in the normal course of investment decision making.

A large body of evidence exists which indicates that security prices incorporate in an unbiased manner all publicly available information and much of what might be called "private information".[69] There is also a large body of evidence which indicates that the security analysis activities of mutual funds and other institutional investors are not reflected in portfolio returns, i.e., they do not increase risk adjusted portfolio returns over a naive random selection buy and hold strategy.[70] Therefore some have been tempted to conclude that the resources expended on such research activities to find under- or over-valued securities is a social loss. Jensen (1974) argues that this conclusion cannot be

[69]See Fama (1970) for a survey of this 'efficient markets' literature.
[70]See Jensen (1969) for an example of this evidence and references.

unambiguously drawn because there is a large consumption element in the demand for these services.

Furthermore, the analysis of this paper would seem to indicate that to the extent that security analysis activities reduce the agency costs associated with the separation of ownership and control they are indeed socially productive. Moreover, if this is true we expect the major benefits of the security analysis activity to be reflected in the higher capitalized value of the ownership claims to corporations and *not* in the period to period portfolio returns of the analyst. Equilibrium in the security analysis industry requires that the private returns to analysis (i.e., portfolio returns) must be just equal to the private costs of such activity,[71] and this will not reflect the social product of this activity which will consist of larger output and higher *levels* of the capital value of ownership claims. Therefore, the argument implies that if there is a non-optimal amount of security analysis being performed it is too much[72] not too little (since the share-holders would be willing to pay directly to have the "optimal" monitoring performed), and we don't seem to observe such payments.

### 6.5. Specialization in the use of debt and equity

Our previous analysis of agency costs suggests at least one other testable hypothesis: i.e., that in those industries where the incentive effects of outside equity or debt are widely different, we would expect to see specialization in the use of the low agency cost financing arrangement. In industries where it is relatively easy for managers to lower the mean value of the outcomes of the enterprise by outright theft, special treatment of favored customers, ease of consumption of leisure on the job, etc. (for example, the bar and restaurant industry) we would expect to see the ownership structure of firms characterized by relatively little outside equity (i.e., 100 percent ownership of the equity by the manager) with almost all outside capital obtained through the use of debt.

The theory predicts the opposite would be true where the incentive effects of debt are large relative to the incentive effects of equity. Firms like conglomerates, in which it would be easy to shift outcome distributions adversely for bond-holders (by changing the acquisition or divestiture policy) should be character-ized by relatively lower utilization of debt. Conversely in industries where the freedom of management to take riskier projects is severely constrained (for example, regulated industries such as public utilities) we should find more intensive use of debt financing.

The analysis suggests that in addition to the fairly well understood role of uncertainty in the determination of the quality of collateral there is at least one other element of great importance – the ability of the owner of the collateral to

---

[71]Ignoring any pure consumption elements in the demand for security analysis.

[72]Again ignoring the value of the pure consumption elements in the demand for security analysis.

change the distribution of outcomes by shifting either the mean outcome or the variance of the outcomes. A study of bank lending policies should reveal these to be important aspects of the contractual practices observed there.

### 6.6. Application of the analysis to the large diffuse ownership corporation

While we believe the structure outlined in the proceeding pages is applicable to a wide range of corporations it is still in an incomplete state. One of the most serious limitation of the analysis is that as it stands we have not worked out in this paper its application to the very large modern corporation whose managers own little or no equity. We believe our approach can be applied to this case but space limitations precludes discussion of these issues here. They remain to be worked out in detail and will be included in a future paper.

### 6.7. The supply side of the incomplete markets question

The analysis of this paper is also relevant to the incomplete market issue considered by Arrow (1964), Diamond (1967), Hakansson (1974a, b), Rubinstein (1974), Ross (1974) and others. The problems addressed in this literature derive from the fact that whenever the available set of financial claims on outcomes in a market fails to span the underlying state space [see Arrow (1964) and Debreu (1959)] the resulting allocation is Pareto inefficient. A disturbing element in this literature surrounds the fact that the inefficiency conclusion is generally drawn without explicit attention in the analysis to the costs of creating new claims or of maintaining the expanded set of markets called for to bring about the welfare improvement.

The demonstration of a possible welfare improvement from the expansion of the set of claims by the introduction of new basic contingent claims or options can be thought of as an analysis of the demand conditions for new markets. Viewed from this perspective, what is missing in the literature on this problem is the formulation of a positive analysis of the supply of markets (or the supply of contingent claims). That is, what is it in the maximizing behavior of individuals in the economy that causes them to create and sell contingent claims of various sorts?

The analysis in this paper can be viewed as a small first step in the direction of formulating an analysis of the supply of markets issue which is founded in the self-interested maximizing behavior of individuals. We have shown why it is in the interest of a wealth maximizing entrepreneur to create and sell claims such as debt and equity. Furthermore, as we have indicated above, it appears that extensions of these arguments will lead to a theory of the supply of warrants, convertible bonds and convertible preferred stock. We are not suggesting that the specific analysis offered above is likely to be sufficient to lead to a theory of the supply of the wide range of contracts (both existing and merely potential) in

the world at large. However, we do believe that framing the question of the completeness of markets in terms of the joining of both the demand and supply conditions will be very fruitful instead of implicitly assuming that new claims spring forth from some (costless) well head of creativity unaided or unsupported by human effort.

## 7. Conclusions

The publicly held business corporation is an awesome social invention. Millions of individuals voluntarily entrust billions of dollars, francs, pesos, etc., of personal wealth to the care of managers on the basis of a complex set of contracting relationships which delineate the rights of the parties involved. The growth in the use of the corporate form as well as the growth in market value of established corporations suggests that at least, up to the present, creditors and investors have by and large not been disappointed with the results, despite the agency costs inherent in the corporate form.

Agency costs are as real as any other costs. The level of agency costs depends among other things on statutory and common law and human ingenuity in devising contracts. Both the law and the sophistication of contracts relevant to the modern corporation are the products of a historical process in which there were strong incentives for individuals to minimize agency costs. Moreover, there were alternative organizational forms available, and opportunities to invent new ones. Whatever its shortcomings, the corporation has thus far survived the market test against potential alternatives.

## References

Alchian, A.A., 1965, The basis of some recent advances in the theory of management of the firm, Journal of Industrial Economics, Nov., 30–44.

Alchian, A.A., 1968, Corporate management and property rights, in: Economic policy and the regulation of securities (American Enterprise Institute, Washington, DC).

Alchian, A.A., 1974, Some implications of recognition of property right transactions costs, unpublished paper presented at the First Interlaken Conference on Analysis and Ideology, June.

Alchian, A.A. and W.R. Allen, 1969, Exchange and production: Theory in use (Wadsworth, Belmont, CA).

Alchian, A.A. and H. Demsetz, 1972, Production, information costs, and economic organization, American Economic Review LXII, no. 5, 777–795.

Alchian, A.A. and R.A. Kessel, 1962, Competition, monopoly and the pursuit of pecuniary gain, in: Aspects of labor economics (National Bureau of Economic Research, Princeton, NJ).

Arrow, K.J., 1963/4, Control in large organizations, Management Science 10, 397–408.

Arrow, K.J., 1964, The role of securities in the optimal allocation of risk bearing, Review of Economic studies 31, no. 86, 91–96.

Atkinson, T.R., 1967, Trends in corporate bond quality, in: Studies in corporate bond finance 4 (National Bureau of Economic Research, New York).

Baumol, W.J., 1959, Business behavior, value and growth (Macmillan, New York).

Becker, G., 1957, The economics of discrimination (University of Chicago Press, Chicago, IL).

Becker, G.S. and G.J. Stigler, 1972, Law enforcement, corruption and compensation of enforcers, unpublished paper presented at the Conference on Capitalism and Freedom, Oct.

Benston, G., 1977, The impact of maturity regulation on high interest rate lenders and borrowers, Journal of Financial Economics 4, no. 1.

Berhold, M., 1971, A theory of linear profit sharing incentives, Quarterly Journal of Economics LXXXV, Aug., 460–482.

Berle, A.A., Jr. and G.C. Means, 1932, The modern corporation and private property (Macmillan, New York).

Black, F. and M. Scholes, 1973, The pricing of options and corporate liabilities, Journal of Political Economy 81, no. 3, 637–654.

Black, F., M.H. Miller and R.A. Posner, 1974, An approach to the regulation of bank holding companies, unpublished manuscript (University of Chicago, Chicago, IL).

Branch, B., 1973, Corporate objectives and market performance, Financial Management, Summer, 24–29.

Coase, R.H., 1937, The nature of the firm, Economica, New Series, IV, 386–405. Reprinted in: Readings in price theory (Irwin, Homewood, IL) 331–351.

Coase, R.H., 1959, The Federal Communications Commission, Journal of Law and Economics II, Oct., 1–40.

Coase, R.H., 1960, The problem of social cost, Journal of Law and Economics III, Oct., 1–44.

Coase, R.H., 1964, Discussion, American Economic Review LIV, no. 3, 194–197.

Cyert, R.M. and C.L. Hedrick, 1972, Theory of the firm: Past, present and future; An interpretation, Journal of Economic Literature X, June, 398–412.

Cyert, R.M. and J.G. March, 1963, A behavioral theory of the firm (Prentice Hall, Englewood Cliffs, NJ).

De Alessi, L., 1973, Private property and dispersion of ownership in large corporations, Journal of Finance, Sept., 839–851.

Debreu, G., 1959, Theory of value (Wiley, New York).

Demsetz, H., 1967, Toward a theory of property rights, American Economic Review LVII, May, 347–359.

Demsetz, H., 1969, Information and efficiency: Another viewpoint, Journal of Law and Economics XII, April, 1–22.

Diamond, P.A., 1967, The role of a stock market in a general equilibrium model with technological uncertainty, American Economic Review LVII, Sept., 759–776.

Evans, J.L. and S.H. Archer, 1968, Diversification and the reduction of dispersion: An empirical analysis, Journal of Finance, Dec.

Fama, E.F., 1970a, Efficient capital markets: A review of theory and empirical work, Journal of Finance XXV, no. 2.

Fama, E.F., 1970b, Multiperiod consumption-investment decisions, American Economic Review LX, March.

Fama, E.F., 1972, Ordinal and measurable utility, in: M.C. Jensen, ed., Studies in the theory of capital markets (Praeger, New York).

Fama, E.F. and M. Miller, 1972, The theory of finance (Holt, Rinehart and Winston, New York).

Friedman, M., 1970, The social responsibility of business is to increase its profits, New York Times Magazine, 13 Sept, 32ff.

Furubotn, E.G. and S. Pejovich, 1972, Property rights and economic theory: A survey of recent literature, Journal of Economic Literature X, Dec., 1137–1162.

Galai, D. and R.W. Masulis, 1976, The option pricing model and the risk factor of stock, Journal of Financial Economics 3, no. 1/2, 53–82.

Hakansson, N.H., 1974a, The superfund: Efficient paths toward a complete financial market, unpublished manuscript.

Hakansson, N.H., 1974b, Ordering markets and the capital structures of firms with illustrations, Institute of Business and Economic Research Working Paper no. 24 (University of California, Berkeley, CA).

Heckerman, D.G., 1975, Motivating managers to make investment decisions, Journal of Financial Economics 2, no. 3, 273–292.

Hirshleifer, J., 1958, On the theory of optimal investment decisions, Journal of Political Economy, Aug., 329–352.

Hirshleifer, J., 1970, Investment, interest, and capital (Prentice-Hall, Englewood Cliffs, NJ).

Jensen, M.C., 1969, Risk, the pricing of capital assets, and the evaluation of investment portfolios, Journal of Business 42, no. 2, 167–247.

Jensen, M.C., 1974, Tests of capital market theory and implications of the evidence. Graduate School of Management Working Paper Series no.7414 (University of Rochester, Rochester, NY).

Jensen, M.C. and J.B. Long, 1972, Corporate investment under uncertainty and Pareto optimality in the capital markets, Bell Journal of Economics, Spring, 151–174.

Jensen, M.C. and W.H. Meckling, 1976, Can the corporation survive? Center for Research in Government Policy and Business Working Paper no. PPS 76–4 (University of Rochester, Rochester, NY).

Klein, W.A., 1976, Legal and economic perspectives on the firm, unpublished manuscript (University of California, Los Angeles, CA).

Kraus, A. and R. Litzenberger, 1973, A state preference model of optimal financial leverage, Journal of Finance, Sept.

Larner, R.J., 1970, Management control and the large corporation (Dunellen, New York).

Lintner, J., 1965, Security prices, risk, and maximal gains from diversification, Journal of Finance XX, Dec., 587–616.

Lloyd-Davies, P., 1975, Risk and optimal leverage, unpublished manuscript (University of Rochester, Rochester, NY).

Long, J.B., 1972, Wealth, welfare, and the price of risk, Journal of Finance, May, 419–433.

Long, J.B., Jr., 1974, Discussion, Journal of Finance XXXIX, no. 12, 485–488.

Machlup, F., 1967, Theories of the firm: Marginalist, behavioral, managerial, American Economic Review, March, 1–33.

Manne, H.G., 1962, The 'higher criticism' of the modern corporation, Columbia Law Review 62, March, 399–432.

Manne, H.G., 1965, Mergers and the market for corporate control, Journal of Political Economy, April, 110–120.

Manne, H.G., 1967, Our two corporate systems: Law and economics, Virginia Law Review 53, March, 259–284.

Manne, H.G., 1972, The social responsibility of regulated utilities, Wisconsin Law Review V, no. 4, 995–1009.

Marris, R., 1964, The economic theory of managerial capitalism (Free Press of Glencoe, Glencoe, IL).

Mason, E.S., 1959, The corporation in modern society (Harvard University Press, Cambridge, MA).

McManus, J.C., 1975, The costs of alternative economic organizations, Canadian Journal of Economics VIII, Aug., 334–350.

Meckling, W.H., 1976, Values and the choice of the model of the individual in the social sciences, Schweizerische Zeitschrift für Volkswirtschaft und Statistik, Dec.

Merton, R.C., 1973, The theory of rational option pricing, Bell Journal of Economics and Management Science 4, no. 1, 141–183.

Merton, R.C., 1974, On the pricing of corporate debt: The risk structure of interest rates, Journal of Finance XXIX, no. 2, 449–470.

Merton, R.C. and M.G. Subrahmanyam, 1974, The optimality of a competitive stock market, Bell Journal of Economics and Management Science, Spring, 145–170.

Miller, M.H. and F. Modigliani, 1966, Some estimates of the cost of capital to the electric utility industry, 1954–57, American Economic Review, June, 333–391.

Modigliani, F. and M.H. Miller, 1958, The costs of capital, corporation finance, and the theory of investment, American Economic Review 48, June, 261–297.

Modigliani, F. and M.H. Miller, 1963, Corporate income taxes and the cost of capital: A correction, American Economic Review June, 433–443.

Monsen, R.J. and A. Downs, 1965, A theory of large managerial firms, Journal of Political Economy, June, 221–236.

Myers, S.C., 1975, A note on the determinants of corporate debt capacity, unpublished manuscript (London Graduate School of Business Studies, London).

Penrose, E., 1958, The theory of the growth of the firm (Wiley, New York).

Preston, L.E., 1975, Corporation and society: The search for a paradigm, Journal of Economic Literature XIII, June, 434–453.

Ross, S.A., 1973, The economic theory of agency: The principals problems, American Economic Review LXII, May, 134–139.

Ross, S.A., 1974a, The economic theory of agency and the principle of similarity, in: M.D. Balch et al., eds., Essays on economic behavior under uncertainty (North-Holland, Amsterdam).

Ross, S.A., 1974b, Options and efficiency, Rodney L. White Center for Financial Research Working Paper no. 3–74 (University of Pennsylvania, Philadelphia, PA).

Rubinstein, M., 1974, A discrete-time synthesis of financial theory, Parts I and II, Institute of Business and Economic Research Working Papers nos. 20 and 21 (University of California, Berkeley, CA).

Scitovsky, T., 1943, A note on profit maximisation and its implications, Review of Economic Studies XI, 57–60.

Sharpe, W.F., 1964, Capital asset prices : A theory of market equilibrium under conditions of risk, Journal of Finance XIX, Sept., 425–442.

Shubik, M., 1970, A curmudgeon's guide to microeconomics, Journal of Economic Literature VIII, June, 405–434.

Silver, M. and R. Auster, 1969, Entrepreneurship, profit and limits on firm size, Journal of Business 42, July, 277–281.

Simon, H.A., 1955, A behavioral model of rational choice, Quarterly Journal of Economics 69, 99–118.

Simon, H.A., 1959, Theories of decision making in economics and behavioral science, American Economic Review, June, 253–283.

Smith, A., 1937, The wealth of nations, Cannan edition (Modern Library, New York).

Smith, C., 1976, Option pricing: A review, Journal of Financial Economics 3, nos. 1/2, 3–52.

Warner, J.B., 1975, Bankruptcy costs, absolute priority, and the pricing of risky debt claims, unpublished manuscript (University of Chicago, Chicago, IL).

Williamson, O.E., 1964, The economics of discretionary behavior: Managerial objectives in a theory of the firm (Prentice-Hall, Englewood Cliffs, NJ).

Williamson, O.E., 1970, Corporate control and business behavior (Prentice-Hall, Englewood Cliffs, NJ).

Williamson, O.E., 1975, Markets and hierarchies: Analysis and antitrust implications (The Free Press, New York).

Wilson, R., 1968, On the theory of syndicates, Econometrica 36, Jan., 119–132.

Wilson, R., 1969, La decision: Agregation et dynamique des orders de preference, Extrait (Editions du Centre National de la Recherche Scientifique, Paris) 288–307.

# Agency Problems and the Theory of the Firm

## Eugene F. Fama

*University of Chicago*

This paper attempts to explain how the separation of security own-
ership and control, typical of large corporations, can be an efficient
form of economic organization. We first set aside the presumption
that a corporation has owners in any meaningful sense. The entre-
preneur is also laid to rest, at least for the purposes of the large
modern corporation. The two functions usually attributed to the
entrepreneur—management and risk bearing—are treated as natu-
rally separate factors within the set of contracts called a firm. The firm
is disciplined by competition from other firms, which forces the
evolution of devices for efficiently monitoring the performance of
the entire team and of its individual members. Individual partici-
pants in the firm, and in particular its managers, face both the
discipline and opportunities provided by the markets for their ser-
vices, both within and outside the firm.

Economists have long been concerned with the incentive problems
that arise when decision making in a firm is the province of managers
who are not the firm's security holders.[1] One outcome has been the
development of "behavioral" and "managerial" theories of the firm
which reject the classical model of an entrepreneur, or owner-

[1] Jensen and Meckling (1976) quote from Adam Smith (1776). The modern literature
on the problem dates back at least to Berle and Means (1932).

manager, who single-mindedly operates the firm to maximize profits, in favor of theories that focus more on the motivations of a manager who controls but does not own and who has little resemblance to the classical "economic man." Examples of this approach are Baumol (1959), Simon (1959), Cyert and March (1963), and Williamson (1964).

More recently the literature has moved toward theories that reject the classical model of the firm but assume classical forms of economic behavior on the part of agents within the firm. The firm is viewed as a set of contracts among factors of production, with each factor motivated by its self-interest. Because of its emphasis on the importance of rights in the organization established by contracts, this literature is characterized under the rubric "property rights." Alchian and Demsetz (1972) and Jensen and Meckling (1976) are the best examples. The antecedents of their work are in Coase (1937, 1960).

The striking insight of Alchian and Demsetz (1972) and Jensen and Meckling (1976) is in viewing the firm as a set of contracts among factors of production. In effect, the firm is viewed as a team whose members act from self-interest but realize that their destinies depend to some extent on the survival of the team in its competition with other teams. This insight, however, is not carried far enough. In the classical theory, the agent who personifies the firm is the entrepreneur who is taken to be both manager and residual risk bearer. Although his title sometimes changes—for example, Alchian and Demsetz call him "the employer"—the entrepreneur continues to play a central role in the firm of the property-rights literature. As a consequence, this literature fails to explain the large modern corporation in which control of the firm is in the hands of managers who are more or less separate from the firm's security holders.

The main thesis of this paper is that separation of security ownership and control can be explained as an efficient form of economic organization within the "set of contracts" perspective. We first set aside the typical presumption that a corporation has owners in any meaningful sense. The attractive concept of the entrepreneur is also laid to rest, at least for the purposes of the large modern corporation. Instead, the two functions usually attributed to the entrepreneur, management and risk bearing, are treated as naturally separate factors within the set of contracts called a firm. The firm is disciplined by competition from other firms, which forces the evolution of devices for efficiently monitoring the performance of the entire team and of its individual members. In addition, individual participants in the firm, and in particular its managers, face both the discipline and opportunities provided by the markets for their services, both within and outside of the firm.

**The Irrelevance of the Concept of Ownership of the Firm**

To set a framework for the analysis, let us first describe roles for management and risk bearing in the set of contracts called a firm. Management is a type of labor but with a special role—coordinating the activities of inputs and carrying out the contracts agreed among inputs, all of which can be characterized as "decision making." To explain the role of the risk bearers, assume for the moment that the firm rents all other factors of production and that rental contracts are negotiated at the beginning of each production period with payoffs at the end of the period. The risk bearers then contract to accept the uncertain and possibly negative difference between total revenues and costs at the end of each production period.

When other factors of production are paid at the end of each period, it is not necessary for the risk bearers to invest anything in the firm at the beginning of the period. Most commonly, however, the risk bearers guarantee performance of their contracts by putting up wealth ex ante, with this front money used to purchase capital and perhaps also the technology that the firm uses in its production activities. In this way the risk bearing function is combined with ownership of capital and technology. We also commonly observe that the joint functions of risk bearing and ownership of capital are re-packaged and sold in different proportions to different groups of investors. For example, when front money is raised by issuing both bonds and common stock, the bonds involve a combination of risk bearing and ownership of capital with a low amount of risk bearing relative to the combination of risk bearing and ownership of capital inherent in the common stock. Unless the bonds are risk free, the risk bearing function is in part borne by the bondholders, and ownership of capital is shared by bondholders and stockholders.

However, ownership of capital should not be confused with owner-ship of the firm. Each factor in a firm is owned by somebody. The firm is just the set of contracts covering the way inputs are joined to create outputs and the way receipts from outputs are shared among inputs. In this "nexus of contracts" perspective, ownership of the firm is an irrelevant concept. Dispelling the tenacious notion that a firm is owned by its security holders is important because it is a first step toward understanding that control over a firm's decisions is not neces-sarily the province of security holders. The second step is setting aside the equally tenacious role in the firm usually attributed to the entre-preneur.

**Management and Risk Bearing: A Closer Look**

The entrepreneur (manager–risk bearer) is central in both the Jensen-Meckling and Alchian-Demsetz analyses of the firm. For

example, Alchian-Demsetz state: "The essence of the classical firm is identified here as a contractual structure with: 1) joint input production; 2) several input owners; 3) one party who is common to all the contracts of the joint inputs; 4) who has the right to renegotiate any input's contract independently of contracts with other input owners; 5) who holds the residual claim; and 6) who has the right to sell his central contractual residual status. The central agent is called the firm's owner and the employer" (1972, p. 794).

To understand the modern corporation, it is better to separate the manager, the agents of points 3 and 4 of the Alchian-Demsetz definition of the firm, from the risk bearer described in points 5 and 6. The rationale for separating these functions is not just that the end result is more descriptive of the corporation, a point recognized in both the Alchian-Demsetz and Jensen-Meckling papers. The major loss in retaining the concept of the entrepreneur is that one is prevented from developing a perspective on management and risk bearing as separate factors of production, each faced with a market for its services that provides alternative opportunities and, in the case of management, motivation toward performance.

Thus, any given set of contracts, a particular firm, is in competition with other firms, which are likewise teams of cooperating factors of production. If there is a part of the team that has a special interest in its viability, it is not obviously the risk bearers. It is true that if the team does not prove viable factors like labor and management are protected by markets in which rights to their future services can be sold or rented to other teams. The risk bearers, as residual claimants, also seem to suffer the most direct consequences from the failings of the team. However, the risk bearers in the modern corporation also have markets for their services—capital markets—which allow them to shift among teams with relatively low transaction costs and to hedge against the failings of any given team by diversifying their holdings across teams.

Indeed, portfolio theory tells us that the optimal portfolio for any investor is likely to be diversified across the securities of many firms.[2] Since he holds the securities of many firms precisely to avoid having his wealth depend too much on any one firm, an individual security holder generally has no special interest in personally overseeing the detailed activities of any firm. In short, efficient allocation of risk bearing seems to imply a large degree of separation of security ownership from control of a firm.

On the other hand, the managers of a firm rent a substantial lump of wealth—their human capital—to the firm, and the rental rates for

[2] Detailed discussions of portfolio models can be found in Fama and Miller (1972, chaps. 6 and 7), Jensen (1972), and Fama (1976, chaps. 7 and 8).

their human capital signaled by the managerial labor market are likely to depend on the success or failure of the firm. The function of management is to oversee the contracts among factors and to ensure the viability of the firm. For the purposes of the managerial labor market, the previous associations of a manager with success and failure are information about his talents. The manager of a firm, like the coach of any team, may not suffer any immediate gain or loss in current wages from the current performance of his team, but the success or failure of the team impacts his future wages, and this gives the manager a stake in the success of the team.

The firm's security holders provide important but indirect assistance to the managerial labor market in its task of valuing the firm's management. A security holder wants to purchase securities with confidence that the prices paid reflect the risks he is taking and that the securities will be priced in the future to allow him to reap the rewards (or punishments) of his risk bearing. Thus, although an individual security holder may not have a strong interest in directly overseeing the management of a particular firm, he has a strong interest in the existence of a capital market which efficiently prices the firm's securities. The signals provided by an efficient capital market about the values of a firm's securities are likely to be important for the managerial labor market's revaluations of the firm's management.

We come now to the central question. To what extent can the signals provided by the managerial labor market and the capital market, perhaps along with other market-induced mechanisms, discipline managers? We first discuss, still in general terms, the types of discipline imposed by managerial labor markets, both within and outside of the firm. We then analyze specific conditions under which this discipline is sufficient to resolve potential incentive problems that might be associated with the separation of security ownership and control.

### The Viability of Separation of Security Ownership and Control of the Firm: General Comments

The outside managerial labor market exerts many direct pressures on the firm to sort and compensate managers according to performance. One form of pressure comes from the fact that an ongoing firm is always in the market for new managers. Potential new managers are concerned with the mechanics by which their performance will be judged, and they seek information about the responsiveness of the system in rewarding performance. Moreover, given a competitive managerial labor market, when the firm's reward system is not responsive to performance the firm loses managers, and the best are the first to leave.

There is also much internal monitoring of managers by managers themselves. Part of the talent of a manager is his ability to elicit and measure the productivity of lower managers, so there is a natural process of monitoring from higher to lower levels of management. Less well appreciated, however, is the monitoring that takes place from bottom to top. Lower managers perceive that they can gain by stepping over shirking or less competent managers above them. Moreover, in the team or nexus of contracts view of the firm, each manager is concerned with the performance of managers above and below him since his marginal product is likely to be a positive function of theirs. Finally, although higher managers are affected more than lower managers, all managers realize that the managerial labor market uses the performance of the firm to determine each manager's outside opportunity wage. In short, each manager has a stake in the performance of the managers above and below him and, as a consequence, undertakes some amount of monitoring in both directions.

All managers below the very top level have an interest in seeing that the top managers choose policies for the firm which provide the most positive signals to the managerial labor market. But by what mechanism can top management be disciplined? Since the body designated for this function is the board of directors, we can ask how it might be constructed to do its job. A board dominated by security holders does not seem optimal or endowed with good survival properties. Diffuse ownership of securities is beneficial in terms of an optimal allocation of risk bearing, but its consequence is that the firm's security holders are generally too diversified across the securities of many firms to take much direct interest in a particular firm.

If there is competition among the top managers themselves (all want to be the boss of bosses), then perhaps they are the best ones to control the board of directors. They are most directly in the line of fire from lower managers when the markets for securities and managerial labor give poor signals about the performance of the firm. Because of their power over the firm's decisions, their market-determined opportunity wages are also likely to be most affected by market signals about the performance of the firm. If they are also in competition for the top places in the firm, they may be the most informed and responsive critics of the firm's performance.

Having gained control of the board, top management may decide that collusion and expropriation of security holder wealth are better than competition among themselves. The probability of such collusive arrangements might be lowered, and the viability of the board as a market-induced mechanism for low-cost internal transfer of control might be enhanced, by the inclusion of outside directors. The latter might best be regarded as professional referees whose task is to stimulate and oversee the competition among the firm's top mana-

gers. In a state of advanced evolution of the external markets that
buttress the corporate firm, the outside directors are in their turn
disciplined by the market for their services which prices them ac-
cording to their performance as referees. Since such a system of
separation of security ownership from control is consistent with the
pressures applied by the managerial labor market, and since it
likewise operates in the interests of the firm's security holders, it
probably has good survival properties.[3]

This analysis does not imply that boards of directors are likely to be
composed entirely of managers and outside directors. The board is
viewed as a market-induced institution, the ultimate internal monitor
of the set of contracts called a firm, whose most important role is to
scrutinize the highest decision makers within the firm. In the team or
nexus of contracts view of the firm, one cannot rule out the evolution
of boards of directors that contain many different factors of produc-
tion (or their hired representatives), whose common trait is that their
marginal products are affected by those of the top decision makers.
On the other hand, one also cannot conclude that all such factors will
naturally show up on boards since there may be other market-induced
institutions, for example, unions, that more efficiently monitor mana-
gers on behalf of specific factors. All one can say is that in a competi-
tive environment lower-cost sets of monitoring mechanisms are likely
to survive. The role of the board in this framework is to provide a
relatively low-cost mechanism for replacing or reordering top mana-
gers; lower cost, for example, than the mechanism provided by an
outside takeover, although, of course, the existence of an outside
market for control is another force which helps to sensitize the inter-
nal managerial labor market.

The perspective suggested here owes much to, but is nevertheless
different from, existing treatments of the firm in the property rights
literature. Thus, Alchian (1969) and Alchian and Demsetz (1972)
comment insightfully on the disciplining of management that takes
place through the inside and outside markets for managers. However,
they attribute the task of disciplining management primarily to the
risk bearers, the firm's security holders, who are assisted to some
extent by managerial labor markets and by the possibility of outside
takeover. Jensen and Meckling (1976) likewise make control of man-

---

[3] Watts and Zimmerman (1978) provide a similar description of the market-induced
evolution of "independent" outside auditors whose function is to certify and, as a
consequence, stimulate the viability of the set of contracts called the firm. Like the
outside directors, the outside auditors are policed by the market for their services which
prices them in large part on the basis of how well they resist perverting the interests of
one set of factors (e.g., security holders) to the benefit of other factors (e.g., manage-
ment). Like the professional outside director, the welfare of the outside auditor de-
pends largely on "reputation."

agement the province of the firm's risk bearers, but they do not allow for any assistance from the managerial labor market. Of all the authors in the property-rights literature, Manne (1965, 1967) is most concerned with the market for corporate control. He recognizes that with diffuse security ownership management and risk bearing are naturally separate functions. But for him, disciplining management is an "entrepreneurial job" which in the first instance falls on a firm's organizers and later on specialists in the process of outside takeover.

When management and risk bearing are viewed as naturally separate factors of production, looking at the market for risk bearing from the viewpoint of portfolio theory tells us that risk bearers are likely to spread their wealth across many firms and so not be interested in directly controlling the management of any individual firm. Thus, models of the firm, like those of Alchian-Demsetz and Jensen-Meckling, in which the control of management falls primarily on the risk bearers, are not likely to allay the fears of those concerned with the apparent incentive problems created by the separation of security ownership and control. Likewise, Manne's approach, in which the control of management relies primarily on the expensive mechanism of an outside takeover, offers little comfort. The viability of the large corporation with diffuse security ownership is better explained in terms of a model where the primary disciplining of managers comes through managerial labor markets, both within and outside of the firm, with assistance from the panoply of internal and external monitoring devices that evolve to stimulate the ongoing efficiency of the corporate form, and with the market for outside takeovers providing discipline of last resort.

### The Viability of Separation of Security Ownership and Control: Details

The preceding is a general discussion of how pressure from managerial labor markets helps to discipline managers. We now examine somewhat more specifically conditions under which the discipline imposed by managerial labor markets can resolve potential incentive problems associated with the separation of security ownership and control of the firm.

To focus on the problem we are trying to solve, let us first examine the situation where the manager is also the firm's sole security holder, so that there is clearly no incentive problem. When he is sole security holder, a manager consumes on the job, through shirking, perquisites, or incompetence, to the point where these yield marginal expected utility equal to that provided by an additional dollar of wealth usable for consumption or investment outside of the firm. The man-

ager is induced to make this specific decision because he pays directly for consumption on the job; that is, as manager he cannot avoid a full ex post settling up with himself as security holder.

In contrast, when the manager is no longer sole security holder, and in the absence of some form of full ex post settling up for deviations from contract, a manager has an incentive to consume more on the job than is agreed in his contract. The manager perceives that, on an ex post basis, he can beat the game by shirking or consuming more perquisites than previously agreed. This does not necessarily mean that the manager profits at the expense of other factors. Rational managerial labor markets understand any shortcomings of available mechanisms for enforcing ex post settling up. Assessments of ex post deviations from contract will be incorporated into contracts on an ex ante basis; for example, through an adjustment of the manager's wage.

Nevertheless, a game which is fair on an ex ante basis does not induce the same behavior as a game in which there is also ex post settling up. Herein lie the potential losses from separation of security ownership and control of a firm. There are situations where, with less than complete ex post settling up, the manager is induced to consume more on the job than he would like, given that on average he pays for his consumption ex ante.

Three general conditions suffice to make the wage revaluation imposed by the managerial labor market a form of full ex post settling up which resolves the managerial incentive problem described above. The first condition is that a manager's talents and his tastes for consumption on the job are not known with certainty, are likely to change through time, and must be imputed by managerial labor markets at least in part from information about the manager's current and past performance. Since it seems to capture the essence of the task of managerial labor markets in a world of uncertainty, this assumption is no real restriction.

The second assumption is that managerial labor markets appropriately use current and past information to revise future wages and understand any enforcement power inherent in the wage revision process. In short, contrary to much of the literature on separation of security ownership and control, we impute efficiency or rationality in information processing to managerial labor markets. In defense of this assumption, we note that the problem faced by managerial labor markets in revaluing the managers of a firm is much entwined with the problem faced by the capital market in revaluing the firm itself. Although we do not understand all the details of the process, available empirical evidence (e.g., Fama 1976, chaps. 5 and 6) suggests that the capital market generally makes rational assessments of the value of

the firm in the face of imprecise and uncertain information. This does not necessarily mean that information processing in managerial labor markets is equally efficient or rational, but it is a warning against strong presumptions to the contrary.

The final and key condition for full control of managerial behavior through wage changes is that the weight of the wage revision process is sufficient to resolve any potential problems with managerial incentives. In this general form, the condition amounts to assuming the desired result. More substance is provided by specific examples.

### Example 1: Marketable Human Capital

Suppose a manager's human capital, his stream of future wages, is a marketable asset. Suppose the manager perceives that, because of the consequent revaluations of future wages, the current value of his human capital changes by at least the amount of an unbiased assessment of the wealth changes experienced by other factors, primarily the security holders, because of his current deviations from contract. Then, as long as the manager is not a risk preferrer, these revaluations of his human capital are a form of full ex post settling up. The manager need not be charged ex ante for presumed ex post deviations from contract since the weight of the wage revision process is sufficient to neutralize his incentives to deviate.

It is important to consider why the manager might perceive that the value of his human capital changes by at least the amount of an unbiased assessment of the wealth changes experienced by other factors due to his deviations from contract. Note first that the market's assessment of such wealth changes is also its assessment of the difference between the manager's ex post marginal product and the marginal product he contracted to deliver ex ante. However, any assessment of the manager's marginal product is likely to include extraneous noise which has little to do with his talents and efforts. Without specific details on what the market takes to be the statistical process governing the evolution of the manager's talents and his tastes for consumption on the job, one cannot say exactly how far it will go in adjusting his future wages to reflect its most recent measurement of his marginal product. Assuming the market uses information rationally, the adjustment is closer to complete the larger the signal in the most recent measurement relative to the noise, but as long as there is some noise in the process, the adjustment is less than complete.[4]

Although his next wage may not adjust by the full amount of an unbiased assessment of the current cost of his deviations from con-

---

[4] Specific illustrations of this point are provided later.

tract, a manager with a multiperiod horizon may perceive that the implied current wealth change, the present value of likely changes in the stream of future wages, is at least as great as the cost of his deviations from contract. In this case, the contemporaneous change in his wealth implied by an eventual adjustment of future wages is a form of full ex post settling up which results in full enforcement of his contract. Moreover, the wage revision process resolves any potential problems about a manager's incentives even though the implied ex post settling up need not involve the firm currently employing the manager; that is, lower or higher future wages due to current deviations from contract may come from other firms.

Of course, changes in a manager's wealth as a consequence of anticipated future wage revisions are not always equivalent to full ex post settling up. When a manager does not expect to be in the labor market for many future periods, the weight of future wage revisions due to current assessments of performance may amount to substantially less than full ex post settling up. However, it is just as important to recognize that the weight of anticipations about future wages may amount to more than full ex post settling up. There may be situations where the personal wealth change perceived by the manager as a consequence of deviations from contract is greater than the wealth change experienced by other factors. Since many readers have had trouble with this point, it is well to bring it closer to home.

Economists (especially young economists) easily imagine situations where the effects of higher or lower quality of a current article or book on the market value of human capital, through enhancement or lowering of "reputation," are in excess of the effects of quality differences on the market value of the specific work to any publisher. Managers can sometimes have similar perceptions with respect to the implications of current performance for the market value of their human capital.

*Example 2: Stochastic Processes for Marginal Products*

The next example of ex post settling up through the wage revision process is somewhat more formal than that described above. We make specific assumptions about the stochastic evolution of a manager's measured marginal product and about how the managerial labor market uses information from the process to adjust the manager's future wages—in a manner which amounts to precise, full ex post settling up for the results of past performance.

Suppose the manager's measured marginal product for any period $t$ is composed of two terms: (i) an expected value, given his talents, effort exerted during $t$, consumption of perquisites, etc.; and (ii)

89

random noise. The random noise may in part result from measurement error, that is, the sheer difficulty of accurately measuring marginal products when there is team production, but it may also arise in part from the fact that effort exerted and talent do not yield perfectly certain consequences. Moreover, because of the uncertain evolution of the manager's talents and tastes, the expected value of his marginal product is itself a stochastic process. Specifically, we assume that the expected value, $\overline{z}_t$, follows a random walk with steps that are independent of the random noise, $\epsilon_t$, in the manager's measured marginal product, $z_t$. Thus, the measured marginal product,

$$z_t = \overline{z}_t + \epsilon_t, \tag{1}$$

is a random walk plus white noise. For simplicity, we also assume that this process describes the manager's marginal product both in his current employment and in the best alternative employment.

The characteristics (parameters) of the evolution of the manager's marginal product depend to some extent on endogenous variables like effort and perquisites consumed, which are not completely observable. Our purpose is to set up the managerial labor market so that the wage revision process resolves any potential incentive problems that may arise from the endogeneity of $z_t$ in a situation where there is separation of security ownership and control of the firm.

Suppose next that risk bearers are all risk neutral and that 1-period market interest rates are always equal to zero. Suppose also that managerial wage contracts are written so that the manager's wage in any period $t$ is the expected value of his marginal product, $\overline{z}_t$, conditional on past measured values of his marginal product, with the risk bearers accepting the noise $\epsilon_t$, in the ex post measurement of the marginal product. We shall see below that this is an optimal arrangement for our risk-neutral risk bearers. However, it is not necessarily optimal for the manager if he is risk averse. A risk-averse manager may want to sell part of the risk inherent in the uncertain evolution of his expected marginal product to the risk bearers, for example, through a long-term wage contract.

We avoid this issue by assuming that, perhaps because of the more extreme moral hazard problems in long-term contracts (remember that $\overline{z}_t$ is in part under the control of the manager) and the contracting costs to which these moral hazard problems give rise, simple contracts in which the manager's wage is reset at the beginning of each period are dominant, at least for some nontrivial subset of firms and managers.[5] If we could also assume away any remaining moral hazard

---

[5] Institutions like corporations, that are subject to rapid technological change with a large degree of uncertainty about future managerial needs, may find that long-term

(managerial incentive) problems, then with risk-averse managers, risk-neutral risk bearers, and the presumed fixed recontracting period, the contract which specifies ex ante that the manager will be paid the current expected value of his marginal product dominates any contract where the manager also shares the ex post deviation of his measured marginal product from its ex ante expected value (see, e.g., Spence and Zeckhauser 1971).

However, contracts which specify ex ante that the manager will be paid the current expected value of his marginal product seem to leave the typical moral hazard problem that arises when there is less than complete ex post enforcement of contracts. The noise $\epsilon_t$ in the manager's marginal product is borne by the risk bearers. Once the manager's expected marginal product $\bar{z}_t$ (= his current wage) has been assessed, he seems to have an incentive to consume more perquisites and provide less effort than are implied in $\bar{z}_t$.

A mechanism for ex post enforcement is, however, built into the model. With the expected value of the manager's marginal product wandering randomly through time, future assessments of expected marginal products (and thus of wages) will be determined in part by $\epsilon_t$, the deviation of the current measured marginal product from its ex ante expected value. In the present scenario, where $\bar{z}_t$ is assumed to follow a random walk, Muth (1960) has shown that the expected value of the marginal product evolves according to

$$\bar{z}_t = \bar{z}_{t-1} + (1 - \phi)\epsilon_{t-1}, \tag{2}$$

where the parameter $\phi$ $(0 < \phi < 1)$ is closer to zero the smaller the variance of the noise term in the marginal product equation (1) relative to the variance of the steps in the random walk followed by the expected marginal product.

In fact, the process by which future expected marginal products are adjusted on the basis of past deviations of marginal products from their expected values leads to a precise form of full ex post settling up. This is best seen by writing the marginal product $z_t$ in its inverted form, that is, in terms of past marginal products and the current noise. The inverted form for our model, a random walk embedded in random noise, is

$$z_t = (1 - \phi)z_{t-1} + \phi(1 - \phi)z_{t-2} + \phi^2(1 - \phi)z_{t-3} + \ldots + \epsilon_t, \tag{3}$$

managerial contracts can only be negotiated at high cost. On the other hand, institutions like governments, schools, and universities may be able to forecast more reliably their future needs for managers (and other professionals) and so may be able to offer long-term contracts at relatively low cost. These institutions can then be expected to attract the relatively risk-averse members of the professional labor force, while the riskier employment offered by corporations attracts those who are willing to accept shorter-term contracts.

so that

$$\bar{z}_t = (1 - \phi)z_{t-1} + \phi(1 - \phi)z_{t-2} + \phi^2(1 - \phi)z_{t-3} + \ldots \qquad (4)$$

(see, e.g., Nelson 1973, chap. 4, or Muth 1960).

For our purposes, the interesting fact is that, although he is paid his ex ante expected marginal product, the manager does not get to avoid his ex post marginal product. For example, we can infer from (4) that $z_{t-1}$ has weight $1 - \phi$ in $\bar{z}_t$; then it has weight $\phi(1 - \phi)$ in $\bar{z}_{t+1}$, $\phi^2(1 - \phi)$ in $\bar{z}_{t+2}$, and so on. In the end, the sum of the contributions of $z_{t-1}$ to future expected marginal products, and thus to future wages, is exactly $z_{t-1}$. With zero interest rates, this means that the risk bearers simply allow the manager to smooth his marginal product across future periods at the going opportunity cost of all such temporal wealth transfers. As a consequence, the manager has no incentive to try to bury shirking or consumption of perquisites in his ex post measured marginal product.

Since the managerial labor market is presumed to understand the weight of the wage revision process, which in this case amounts to precise full ex post settling up, any potential managerial incentive problems in the separation of risk bearing, or security ownership, from control are resolved. The manager can contract for and take an optimal amount of consumption on the job. The wage set ex ante need not include any allowance for ex post incentives to deviate from the contract since the wage revision process neutralizes any such incentives. Note, moreover, that the value of $\phi$ in the wage revision process described by (4) determines how the observed marginal product of any given period is subdivided and spread across future periods, but whatever the value of $\phi$, the given marginal product is fully accounted for in the stream of future wages. Thus, it is now clear what was meant by the earlier claim that although the parameter $\phi$ in the process generating the manager's marginal product is to some extent under his control, this is not a matter of particular concern to the managerial labor market.

A somewhat evident qualification is in order. The smoothing process described by (4) contains an infinite number of terms, whereas any manager has a finite working life. For practical purposes, full ex post settling up is achieved as long as the manager's current marginal product is "very nearly" fully absorbed by the stream of wages over his future working life. This requires a value of $\phi$ in (4) which is sufficiently far from 1.0, given the number of periods remaining in the manager's working life. Recall that $\phi$ is closer to 1.0 the larger the variance of the noise in the manager's measured marginal product relative to the variance of the steps of the random walk taken by the expected value of his marginal product. Intuitively, when the variance

of the noise term is large relative to that of the changes in the expected value, the current measured marginal product has a weak signal about any change in the expected value of the marginal product, and the current marginal product is only allocated slowly to expected future marginal products.

*Some Extensions*

Having qualified the analysis, let us now indicate some ways in which it is robust to changes in details of the model.

1.  More Complicated Models for the Manager's Marginal Product

The critical ingredient in enforcing precise full ex post settling up through wage revisions on the basis of reassessments of expected marginal products is that when the marginal product and its expected value are expressed in inverted form, as in (3) and (4), the sum of the weights on past marginal products is exactly 1.0. This will be the case (see, e.g., Nelson 1973, chap. 4) whenever the manager's marginal product conforms to a nonstationary stochastic process, but the changes from period to period in the marginal product conform to some stationary ARMA (mixed autoregressive moving average) process. The example summarized in equations (1)–(4) is the interesting but special case where the expected marginal product follows a random walk so that the differences of the marginal product are a stationary, first-order moving average process. The general case allows the expected value of the marginal product to follow any more complicated nonstationary process which has the property that the differences of the marginal product are stationary, so that the marginal product and its expected value can be expressed in inverted form as

$$z_t = \pi_1 z_{t-1} + \pi_2 z_{t-2} + \ldots + \epsilon_t \tag{5}$$

$$\overline{z}_t = \pi_1 z_{t-1} + \pi_2 z_{t-2} + \ldots \tag{6}$$

with

$$\sum_{i=1}^{\infty} \pi_i = 1. \tag{7}$$

These can be viewed as the general conditions for enforcing precise full ex post settling through the wage revision process when the

93

manager's wage is equal to the current expected value of his marginal product.[6]

## 2.  Risk-Averse Risk Bearers

In the framework summarized in equations (5)–(7), if the manager switches firms, the risk bearers of his former firm are left with the remains of his measured marginal products not previously absorbed into the expected value of his marginal product. Nevertheless, in the way we have set up the world, the risk bearers realize that the manager's next firm continues to set his wage according to the same stochastic process as the last firm. Since this results in full ex post settling up on the part of the manager, the motive for switching firms cannot be to avoid perverse adjustments of future wages on the basis of past performance. On average, the switching of managers among firms does not result in gains or losses to risk bearers, which means that the switches are a matter of indifference to our presumed risk-neutral risk bearers.

It is, however, interesting to examine how the analysis might change when the risk bearers are risk averse and switching of managers among firms is not a matter of indifference. Suppose, for the moment, that the risk bearers offer managers contracts where, as before, the manager's wage tracks the expected value of his marginal product, but each period there is also a fixed discount in the wage to compensate the risk bearers for the risks of unfinished ex post settling up with the firm as a consequence of a possible future shift by the manager to another firm. Such an arrangement may satisfy the risk bearers, but it will not be acceptable to the manager. As long as his marginal product evolves according to equations (5)–(7), both in his current firm and in the best alternative, the manager is subject to full ex post settling up. Thus, any risk adjustment of his wage to reflect the fact that the settling up may not be with his current firm is an uncompensated loss which he will endeavor to avoid.

The manager can avoid any risk discount in his wage, and maintain complete freedom to switch among firms, by himself bearing all the risk of his marginal product; that is, he contracts to accept, at the end of each period, his ex post measured marginal product rather than its ex ante expected value so that there is, period by period, full ex post settling up with his current firm. There is such a presumption against

---

[6] When $\bar{z}_t$ follows a stationary process, the long-run average value toward which the process always tends will eventually be known with near perfect certainty. Thus, the case of a stationary expected marginal product is of little interest, at least for the purposes of ex post settling up enforced by the wage revision process.

the optimality of immediate, full ex post settling up in the literature on optimal contracting that it behooves us to examine how and why it works, and is optimal, in the circumstances under examination.

### Contractual Settling Up

The literature on optimal contracting, for example, Harris and Raviv (1978, 1979), Holmström (1979), and Shavell (1979), suggests uniformly that when there is noise in the manager's marginal product, that is, when the deviation of measured marginal product from its expected value cannot be traced unambiguously and costlessly to the manager's actions (talents, effort exerted, and consumption on the job), then a risk-averse manager will always choose to share part of the uncertainty in the evaluation of his performance with the firm's risk bearers. He will agree to some amount of ex post settling up, but always less than 100 percent of the deviation of his measured marginal product from its ex ante expected value. In short, the contracting models suggest that we must learn to live with the incentive problems that arise when there is less than complete ex post enforcement of contracts.

The contracting literature is almost uniformly concerned with 1-period models. In a 1-period world, there can be no enforcement of contracts through a wage revision process imposed by the managerial labor market. The existence of this form of ex post settling up in a multiperiod world affects the manager's willingness to engage in explicit contractual ex post settling up.

For example, in the model summarized in equations (5)–(7), the manager's wage in any period is the expected value of his marginal product assessed at the beginning of the period, and the manager does not immediately share any of the deviation of his ex post marginal product from its ex ante expected value. However, because it contains information about future expected values of his marginal product, eventually the manager's current measured marginal product is allocated in full to future expected marginal products. Equivalently, in the wage revision process described by equations (5)–(7), the managerial labor market in effect acts as a financial intermediary. It withdraws portions of past accumulated measured marginal products to pay the manager a dividend on his human capital equal to the expected value of his marginal product, and implicitly provides the lending arrangements which allow the manager to spread his current measured marginal product over future periods in precisely the way the current marginal product will contribute to expected future marginal products.

Looked at from this perspective, however, the manager might simply contract to take the ex post measured value of his marginal product as his wage and then himself use the capital market to smooth his measured marginal product over future periods. Since the same asset (his human capital) is involved, the manager should be able to carry out these smoothing transactions via the capital market on the same terms as can be had in the managerial labor market. The advantage to the manager in smoothing through the capital market, however, is that he can then contract to accept full ex post settling up period by period (he is paid his measured marginal product), which means he can avoid any risk discount in his wage that might be imposed when he is paid the expected value of his marginal product with the possibility of unanticipated switches to other firms.[7]

It is important to recognize that using the capital market in the manner described above allows the manager to "average out" the random noise in his measured marginal product. Thus, when he is instead paid the expected value of his marginal product each period, and when the process generating his marginal product is described by equations (5)–(7), the manager's current measured marginal product is eventually allocated in full to future expected marginal products. This happily, but only coincidentally, resolves incentive problems by imposing full ex post settling up. The allocation of the current marginal product to future expected marginal products in fact occurs because the current marginal product has information about future expected marginal products. The weights $\pi_i$ in equations (5)–(7) are precisely those that optimally extract this information and so optimally smooth or average out the purely random noise in the manager's measured marginal product. The manager can achieve the same result by contracting to be paid the measured value of his marginal product and then using the capital market to smooth his marginal product. This power of the capital market to reduce the terror in full contractual ex post settling up is lost in the 1-period models that dominate the contracting literature.

---

[7] With positive interest rates, contracting to be paid his measured marginal product and then using the capital market to smooth the marginal product over future periods dominates the contract in which the manager is paid the expected value of his marginal product. Equivalence can be restored by adjusting the expected marginal product $\bar{z}_t$ in eq. (6) for accumulated interest on the past marginal products, $z_{t-1}, z_{t-2}, \ldots$, or by prepaying the present value of interest on the deferrals of the current marginal product over future periods. Suffice it to say, however, that either accumulation or prepayment of interest complicates the problems posed by possible shifts of the manager to other firms and so may lean the system toward contracts in which the manager is paid his measured marginal product and then uses the capital market to achieve optimal smoothing.

## Conclusions

The model summarized by equations (5)–(7) is one specific scenario in which the wage revision process imposed by the managerial labor market amounts to full ex post settling up by the manager for his past performance. The important general point is that in any scenario where the weight of the wage revision process is at least equivalent to full ex post settling up, managerial incentive problems—the problems usually attributed to the separation of security ownership and control of the firm—are resolved.

No claim is made that the wage revision process always results in a full ex post settling up on the part of the manager. There are certainly situations where the weight of anticipated future wage changes is insufficient to counterbalance the gains to be had from ex post shirking, or perhaps outright theft, in excess of what was agreed ex ante in a manager's contract. On the other hand, precise full ex post settling up is not an upper bound on the force of the wage revision process. There are certainly situations where, as a consequence of anticipated future wage changes, a manager perceives that the value of his human capital changes by more than the wealth changes imposed on other factors, and especially the firm's security holders, by his current deviations from the terms of his contract.

The extent to which the wage revision process imposes ex post settling up in any particular situation is, of course, an empirical issue. But it is probably safe to say that the general phenomenon is at least one of the ingredients in the survival of the modern large corporation, characterized by diffuse security ownership and the separation of security ownership and control, as a viable form of economic organization.

### References

Alchian, Armen A. "Corporate Management and Property Rights." In *Economic Policy and the Regulation of Corporate Securities*, edited by Henry G. Manne. Washington: American Enterprise Inst. Public Policy Res., 1969.
Alchian, Armen A., and Demsetz, Harold. "Production, Information Costs, and Economic Organization." *A.E.R.* 62 (December 1972): 777–95.
Baumol, William J. *Business Behavior, Value and Growth.* New York: Macmillan, 1959.
Berle, Adolph A., Jr., and Means, Gardiner C. *The Modern Corporation and Private Property.* New York: Macmillan, 1932.
Coase, Ronald H. "The Nature of the Firm." *Economica*, n.s. 4 (November 1937): 386–405.
———. "The Problem of Social Cost." *J. Law and Econ.* 3 (October 1960): 1–44.
Cyert, Richard M., and March, James G. *A Behavioral Theory of the Firm.* Englewood Cliffs, N.J.: Prentice-Hall, 1963.

Fama, Eugene F. *Foundations of Finance*. New York: Basic, 1976.

Fama, Eugene F., and Miller, Merton H. *The Theory of Finance*. New York: Holt, Rinehart & Winston, 1972.

Harris, Milton, and Raviv, Artur. "Some Results on Incentive Contracts with Applications to Education and Employment, Health Insurance, and Law Enforcement." *A.E.R.* 68 (March 1978): 20–30.

———. "Optimal Incentive Contracts with Imperfect Information." Working Paper no. 70-75-76, Carnegie-Mellon Univ., Graduate School of Indus. Admin., April 1976 (rev. January 1979), forthcoming in *J. Econ. Theory*.

Holmström, Bengt. "Moral Hazard and Observability." *Bell J. Econ.* 10 (Spring 1979): 74–91.

Jensen, Michael C. "Capital Markets: Theory and Evidence." *Bell J. Econ. and Management Sci.* 3 (Autumn 1972): 357–98.

Jensen, Michael C., and Meckling, William H. "Theory of the Firm: Managerial Behavior, Agency Costs and Ownership Structure." *J. Financial Econ.* 3 (October 1976): 305–60.

Manne, Henry G. "Mergers and the Market for Corporate Control." *J.P.E.* 73, no. 2 (April 1965): 110–20.

———. "Our Two Corporate Systems: Law and Economics." *Virginia Law Rev.* 53 (March 1967): 259–85.

Muth, John F. "Optimal Properties of Exponentially Weighted Forecasts." *J. American Statis. Assoc.* 55 (June 1960): 299–306.

Nelson, Charles R. *Applied Time Series Analysis for Managerial Forecasting*. San Francisco: Holden-Day, 1973.

Shavell, Steven. "Risk Sharing and Incentives in the Principal and Agent Relationship." *Bell J. Econ.* 10 (Spring 1979): 55–73.

Simon, Herbert A. "Theories of Decision Making in Economics and Behavioral Science." *A.E.R.* 49 (June 1959): 253–83.

Smith, Adam. *The Wealth of Nations*. 1776. Cannan ed. New York: Modern Library, 1937.

Spence, Michael, and Zeckhauser, Richard. "Insurance, Information and Individual Action." *A.E.R.* 61 (May 1971): 380–87.

Watts, Ross L., and Zimmerman, Jerold. "Auditors and the Determination of Accounting Standards, an Analysis of the Lack of Independence." Working Paper GPB 7806, Univ. Rochester, Graduate School of Management, 1978.

Williamson, Oliver E. *The Economics of Discretionary Behavior: Managerial Objectives in a Theory of the Firm*. Englewood Cliffs, N.J.: Prentice-Hall, 1964.

# AGENCY PROBLEMS
# AND RESIDUAL CLAIMS*

*EUGENE F. FAMA*      and      *MICHAEL C. JENSEN*
*University of Chicago*              *University of Rochester*

## I. Introduction

### A. *Organizational Survival*

Social and economic activities, such as religion, entertainment, education, research, and the production of other goods and services, are carried on by different types of organizations, for example, corporations, proprietorships, partnerships, mutuals, and nonprofits. Most goods and services can be produced by any form of organization, and there is competition among organizational forms for survival in any activity. Absent fiat, the form of organization that survives in an activity is the one that delivers the product demanded by customers at the lowest price while covering costs. This is the telling dimension on which the economic environment chooses among organizational forms.

An important factor in the survival of organizational forms is control of agency problems. Agency problems arise because contracts are not costlessly written and enforced. Agency costs include the costs of structuring, monitoring, and bonding a set of contracts among agents with conflicting interests, plus the residual loss incurred because the cost of full enforcement of contracts exceeds the benefits.[1] In this paper we explain the

---

[1] This definition of agency costs first appears in Michael C. Jensen & William H. Meckling, Theory of the Firm: Managerial Behavior, Agency Costs and Ownership Structure, 3 J. Financial Econ. 305 (1976).

special features of the residual claims of different organizational forms as
efficient approaches to controlling special agency problems. We analyze
only private organizations. In related papers we examine other features of
the contract structures of different organizational forms that contribute to
their survival; in particular, (1) the control of agency problems in the class
of organizations characterized by separation of "ownership" and "con-
trol," and (2) the effects of special characteristics of residual claims on
decision rules for resource allocation.[2]

## B.  Residual Claims: General Discussion

The contract structures of organizations limit the risks undertaken by
most agents by specifying either fixed payoffs or incentive payoffs tied to
specific measures of performance. The residual risk—the risk of the dif-
ference between stochastic inflows of resources and promised payments
to agents—is borne by those who contract for the rights to net cash flows.
We call these agents the residual claimants or residual risk bearers.

The characteristics of residual claims distinguish organizations from
one another and help explain the survival of organizational forms in
specific activities. We first analyze and contrast the relatively unrestricted
residual claims of open corporations with the restricted residual claims of
proprietorships, partnerships, and closed corporations. We then turn to
the more specialized residual claims of professional partnerships,
financial mutuals, and nonprofits.

## II.  OPEN CORPORATIONS

Most large nonfinancial organizations are open corporations. The com-
mon stock residual claims of such organizations are unrestricted in the
sense that (1) stockholders are not required to have any other role in
the organization, (2) their residual claims are freely alienable, and (3) the
residual claims are rights in net cash flows for the life of the organization.
Because of the unrestricted nature of the residual claims of open corpora-
tions, there is generally almost complete separation and specialization of
decision functions and residual risk bearing.

## A.  Common Stock versus State Contingent Claims

One can imagine claims that are even less restricted than the common
stocks of open corporations. There could be "state contingent claims"—

that is, claims of the sort discussed by Arrow and Debreu[3] specifying payoffs for each possible future state of the world. Such state contingent claims allow any (hence generally "less restricted") allocation of risk. They are, nonetheless, fixed payoff promises. To specify the total payoffs to be obtained in all future states, one would need to identify all current and future decisions of an organization through state contingent claim contracts. Given the costs and information requirements this implies, it is not surprising that state contingent claims are not the dominant system for allocating risk.

We can also imagine state contingent claims that are true residual claims. The claim would cover a fraction of the organization's net cash flows in a given state rather than a specified payoff in that state. However, this type of claim generates conflicts among the claim holders of different states because alternative decisions shift payoffs across states and benefit some claim holders at the expense of others. Common stock that represents proportionate claims on the payoffs of all future states eliminates these agency problems, but at the sacrifice of some efficiency in the allocation of risk. Common stock and other common forms of residual claims also avoid most of the costs of defining and verifying states of the world.

### B.   The Advantages of Common Stock Residual Claims

1.   *Unrestricted Risk Sharing among Residual Claimants.*   The common stock of open corporations allows more efficient risk sharing than residual claims that are not separable from decision roles, as, for example, in proprietorships and partnerships where the proprietors and partners are the decision makers and the primary residual claimants. Common stock allows residual risk to be spread across many residual claimants who individually choose the extent to which they bear risk and who can diversify across organizations offering such claims. Other things equal, portfolio theory implies that such unrestricted risk sharing lowers the cost of risk-bearing services.[4]

2.   *Specialized Risk Bearing by Residual Claimants.*   The activities of large open nonfinancial corporations are typically complicated, involving contracts with many factors of production, for example, different types of labor, raw materials, and managers. When there is significant variation through time in the probability of default on these contracts, contracting costs increase. In addition, because the human capital of agents is gener-

---

[3] Kenneth J. Arrow, The Role of Securities in the Optimal Allocation of Risk Bearing, 31 Rev. Econ. Stud. 91 (1964); Gerard Debreu, Theory of Value (1959).

[4] See, for example, Arrow, *supra* note 3; or Eugene F. Fama, Foundations of Finance chs. 7 & 8 (1976).

ally employed in a single organization, risk aversion tends to cause them to charge more for any risk they bear than security holders who can diversify risk across many organizations.[5]

Efficient accommodation of large-scale specialized risk bearing by residual claimants is an advantage of corporate common stock. To bond contractual payments to other agents, the common stockholders put up wealth, which is used to purchase assets. If the wealth required to bond promised payments goes beyond the value of inputs optimally purchased rather than rented, common stock proceeds can be used to purchase liquid assets, for example, the securities of other organizations, that have no function except to bond specialization of risk bearing by residual claimants.

3. *Purchase of Organization-specific Assets.* Klein, Crawford, and Alchian and Jensen and Meckling argue that because of conflicts of interest with outside owners of organization-specific assets—assets that have lower value to other organizations—rental contracts for such assets generate higher agency costs than outright purchase.[6] Common stock, with its capacity for raising wealth from residual claimants, is an efficient vehicle for financing such purchases in activities where using large amounts of organization-specific risky assets is efficient.

4. *Specialization of Management.* In the complicated production and distribution activities of large open corporations, coordinating the activities of agents, recontracting among them, and initiating and implementing resource allocation decisions are specialized tasks which are important to the survival of the organization and largely fall on its managers. However, managerial skills are not necessarily tied to wealth or willingness to bear risk, and incompetent managers who are important residual claimants can be difficult to remove. Thus, ignoring agency problems in the decision process, the survival of a complex organization is enhanced by common stock residual claims that allow specialization of management—in effect, the absence of a classical entrepreneur who is both decision maker and residual risk bearer.

5. *The Market Value Rule for Investment Decisions.* When common stocks are traded without transactions costs in a perfectly competitive capital market, the stockholders agree that resource allocation decisions

[5] See Patricia B. Reagan & Rene M. Stulz, Risk Bearing, Labor Contracts, and Capital Markets, (Working Paper Series No. MERC 82-19 Univ. Rochester Managerial Economics Research Center 1982) for an analysis of risk sharing between internal agents and residual claimants and for references to the related literature.

[6] Benjamin Klein, Robert Crawford, & Armen A. Alchian, Vertical Integration, Appropriable Rents, and the Competitive Contracting Process, 21 J. Law & Econ. 297 (1978); Michael C. Jensen & William H. Meckling, Rights and Production Functions: An Application to Labor-managed Firms and Codetermination, 52 J. Bus. 469 (1979).

should be evaluated according to their contribution to the current market value of their residual claims.[7] The market value rule weighs current against future resources according to the opportunity costs at which resources can be traded across time in the capital market. For example, the market value rule favors expenditures to reduce the current and future costs of delivering products whenever the current market value of the future cost savings is greater than the current expenditure. Product prices can then be lowered while still covering costs.

In contrast, when the horizon of the residual claims is less than the life of the organization, residual claimants assign zero value to cash flows that occur beyond the horizon.[8] Similarly, when residual claims are not freely alienable or separable from other roles in the organization, it is rational for risk bearers to attribute lower current value to uncertain cash flows than is implied by capital market prices for the future resources.[9] As a consequence, ignoring agency problems in the decision process, organizations with common stock residual claims, investing according to the market value rule which is optimal for their residual claimants, will be able to deliver products at lower prices than organizations with restricted residual claims.

## C.   The Agency Problems of Common Stock Residual Claims

The unrestricted nature of the common stock residual claims of open corporations leads to an important agency problem. The decision process is in the hands of professional managers whose interests are not identical to those of residual claimants. This problem of separation of "ownership" and "control"—more precisely, the separation of residual risk bearing from decision functions—has troubled students of open corporations from Adam Smith to Berle and Means and Jensen and Meckling.[10] In "Separation of Ownership and Control"[11] we argue that this agency problem is controlled by decision systems that separate the management (initiation and implementation) and control (ratification and monitoring) of important decisions at all levels of the organization.

[7] See, for example, Eugene F. Fama, The Effects of a Firm's Investment and Financing Decisions on the Welfare of its Security Holders, 68 Am. Econ. Rev. 272 (1978).

[8] See E. G. Furubotn & S. Pejovich, Property Rights, Economic Decentralization and the Evolution of the Yugoslav Firm, 1965–1972, 16 J. Law & Econ. 275 (1973); and Jensen & Meckling, *supra* note 8.

[9] The details of the argument are in Fama & Jensen, Organizational Forms, *supra* note 2.

[10] Adam Smith, The Wealth of Nations (Cannan ed. 1904) (1st ed. London 1776); Adolf A. Berle & Gardiner C. Means, The Modern Corporation and Private Property (1932); Jensen & Meckling, *supra* note 1.

[11] Fama & Jensen, in this issue.

Devices for separating decision management and decision control include (1) decision hierarchies in which the decision initiatives of lower level agents are passed on to higher level agents, first for ratification and then for monitoring, (2) boards of directors that ratify and monitor the organization's most important decisions and hire, fire, and compensate top-level decision managers, and (3) incentive structures that encourage mutual monitoring among decision agents. The costs of such mechanisms for separating decision management from decision control are part of the price that open corporations pay for the benefits of unrestricted common stock residual claims.

## III.   RESTRICTED VERSUS UNRESTRICTED RESIDUAL CLAIMS

The proprietorships, partnerships, and closed corporations observed in small-scale production activities differ in many ways both from one another and from open corporations. For example, proprietorships have a single residual claimant, whereas partnerships and closed corporations have multiple residual claimants. As a consequence, the residual claim contracts in partnerships and closed corporations must specify rights in net cash flows and procedures for transferring residual claims to new agents more explicitly than the residual claims in proprietorships.

However, for control of the agency problems in the decision process, the common characteristic of the residual claims of proprietorships, partnerships, and closed corporations that distinguishes them from open corporations is that the residual claims are largely restricted to important decision agents. This restriction avoids the agency problems between residual claimants and decision agents that arise because of separation of risk-bearing and decision functions in open corporations. Thus, costly mechanisms for separating the management and control of decisions are avoided.[12]

Restricting residual claims to decision makers controls agency problems between residual claimants and decision agents, but at the expense of the benefits of unrestricted common stock. The decision process suffers efficiency losses because decision agents must be chosen on the basis of wealth and willingness to bear risk as well as for decision skills. Residual claimants forgo optimal diversification so that residual claims and decision making can be combined in a small number of agents. Forgone diversification and limited alienability lower the value of the residual claims, raise the cost of risk-bearing services, and lead to less investment

---

[12] However, in partnerships and closed corporations, some mechanisms for resolving conflicts among residual claimant decision makers (for example, buy-out rules) are required.

in projects with uncertain payoffs than when residual claims are unrestricted. Finally, because decision agents have limited wealth, restricting residual claims to them also limits resources available for bonding contractual payoffs and for acquiring risky organization-specific assets.

An organizational form survives in an activity when the costs and benefits of its residual claims and the approaches it provides to controlling agency problems combine with available production technology to allow the organization to deliver products at lower prices than other organizational forms. The restricted residual claims of proprietorships, partnerships, and closed corporations are more likely to dominate when technology does not involve important economies of scale that lead to large demands for specialized decision skills, specialized risk bearing, and wealth from residual claimants. In these circumstances, the agency costs saved by restricting residual claims to decision agents outweigh the benefits that would be obtained from separation and specialization of decision and risk-bearing functions. On the other hand, unrestricted common stock residual claims are more likely to dominate when there are important economies of scale in production that (i) can be realized only with a complex decision hierarchy that makes use of specialized decision skills throughout the organization, (ii) generate large aggregate risks to be borne by residual claimants, and (iii) demand large amounts of wealth from residual claimants to purchase risky assets and to bond the payoffs promised to a wide range of agents in the organization. In such complex organizations the benefits of unrestricted common stock residual claims are likely to outweigh the costs of controlling the agency problems inherent in the separation and specialization of decision and risk-bearing functions. In these circumstances, the open corporation is more likely to win the competition for survival.[13]

## IV.   Special Forms of Residual Claims

The restriction of residual claims to important decision agents distinguishes the residual claims of proprietorships, partnerships, and closed corporations from the unrestricted residual claims of open corporations. There are, however, other organizational forms, including professional partnerships, financial mutuals, and nonprofits, that offer more unusual residual claims. We explain the special characteristics of the residual claims of these organizations as effective devices for controlling special agency problems.

---

[13] In Fama & Jensen, Separation of Ownership, in this issue, we discuss how the diffusion of information among decision agents influences the survival of organizational forms. For simplicity, we have ignored these issues here.

### A. Professional Partnerships

Like the proprietorships, partnerships, and closed corporations discussed above, the residual claims of the professional partnerships observed in law, public accounting, medicine, and business consulting are restricted to important decision agents. However, in professional partnerships, a partner's share in net cash flows is renegotiated periodically, and his rights in net cash flows are often limited to his period of service in the organization. In effect, a professional partner's residual claim is a flexible and inalienable share of net cash flows for a limited horizon. Flexible sharing rules, inalienability, and limited horizons distinguish the residual claims of professional partnerships from those of the proprietorships, partnerships, and closed corporations observed in other activities. Moreover, these special features of professional partnership residual claims are generally retained when these organizations become professional service corporations for tax purposes.

1. *Decentralized Decision Making and Restricted Residual Claims.* In professional partnerships, large and small, individuals or small teams work on cases, audits, and so on. Because of the importance of specific knowledge about particular clients—knowledge that is costly to transfer among agents—it is efficient for the teams in large partnerships to make most decisions locally. Thus, with respect to the services rendered to customers, decision control takes place within teams, where interaction and mutual monitoring are heaviest. At this level, however, decision management (initiation and implementation) and decision control (ratification and monitoring) are not separate. To control the resulting agency problems, the residual claims in professional partnerships are restricted to the professional agents who are the important team members and who have major decision making roles. This is consistent with the hypothesis developed in "Separation of Ownership and Control"[14] that combination of decision management and control functions in one or a few agents leads to restriction of residual claims to the important decision agents.

2. *The Demand for Monitoring, Bonding, and Consulting.* Lawyers, public accountants, physicians, and some business consultants provide services where one incompetent act can do large damage to a client. As a consequence, certification and pedigree are important to clients. Moreover, even in the largest professional service organizations, services are rendered in individual cases by one or a few professionals. Responsibility for variation in the quality of services is easily assigned to individual agents, and the performance of agents is often well known to clients. In these circumstances, the value of human capital is sensitive to perfor-

[14] *Id.*

mance. In effect, unlimited liability is imposed on the human capital of professional agents by the market for their services. This gives the professional incentives to purchase monitoring and consulting to help limit losses in the value of human capital.

Since professional services are technical, a lawyer, physician, public accountant, or business consultant is efficiently monitored by others of the same training who can also provide valuable consulting services. Such mutual monitoring and consulting are encouraged when professional agents agree to pool net cash flows and to share liability for the actions of colleagues. Pooling of net cash flows and liability is attractive because it encourages mutual monitoring and consulting. Mutual monitoring and consulting improve the quality of services delivered, control liability losses, and enhance the human capital of the partners. Pooling of net cash flows and liability also has risk-sharing advantages.

The analysis is robust to the fact that partnerships sometimes purchase malpractice insurance. Insurance eliminates variability of liability payoffs by substituting a certain insurance premium. However, if premiums are renegotiated to reflect the malpractice experience of the insured, insurance does not destroy the professional's incentives to be monitored or to consult with other professionals.[15] In addition, insurance covers liability to customers but not reductions in the value of human capital caused by incompetent or malfeasant acts.

3. *Large Professional Partnerships and Flexible Sharing Rules.* Some professional partnerships have hundreds and sometimes thousands of partners. Such large partnerships provide portfolios of specialized services that are marketed and delivered over a wide geographical area. They can also provide large bonds to protect clients against losses from malfeasance or incompetence.[16] Large partnerships are also educational organizations, offering young professionals a wide range of opportunities and interaction with other professionals. We are more concerned, though, with the effects of size on the contract structures of these organizations than with explaining why they are large.

Having attained partner status, a professional may be tempted to free-ride on the efforts of colleagues. The residual claims of large partnerships take a direct approach to this agency problem. The residual claim is not generally a fixed share of net cash flows. Rather, a partner's share is renegotiated annually on the basis of past performance and estimates of likely contributions to future net cash flows. In these large partnerships

---

[15] David Mayers & Clifford W. Smith, Jr., On the Corporate Demand for Insurance, 55 J. Bus. 281 (1982), argue that insurance itself is a way to purchase monitoring.

[16] See Linda DeAngelo, Auditor Size and Audit Quality, 3 J. Accounting & Econ. 183 (1981).

service to a client is delivered by a small group of professionals who interact and monitor one another intensively. The composition of the teams changes from case to case to match specialized talents to specialized problems. As a result, the professionals develop knowledge of the talents and contributions of a range of colleagues. Flexible sharing rules add to partners' incentives to gather and communicate such knowledge to the renegotiation process.

Given flexible sharing rules and the way payoffs are tied to performance, large professional partnerships can be viewed as associations of proprietors who get together to obtain the benefits from marketing a portfolio of specialized skills both to clients and to young professionals who purchase specialized education. Or, since the partners often work in small teams that shift from case to case, a large partnership can be regarded as a fluid association of small partnerships.

4. *Limited Horizon Residual Claims.* Limitations on the horizon covered by residual claims cause organizations to bias decisions against alternatives that generate net cash flows beyond the horizon. In "Organizational Forms and Investment Decisions"[17] we argue that the limited horizon feature of the residual claims of professional partnerships reflects the relative unimportance of assets that are not effectively capitalized in the human capital of existing partners. There are generally no important patents, specialized assets, or technologies to be passed from one generation of partners to the next. Each partner brings a depleting asset—human capital—to the partnership. The annual readjustments of shares in net cash flows that are typical, especially in large professional partnerships, calibrate a partner's payoffs to reflect the current and expected future contributions of his human capital. When a partner's human capital is used up or withdrawn from the organization, contributions to net cash flows cease, and this is reflected in the termination, without substantial compensation, of his residual claim.

This explanation of the limited horizon feature of the residual claims of professional partnerships gets support from several sources:

1. Professional human capital serves as a bond against malfeasance when its value is sensitive to performance. However, professional human capital cannot be sold to cover liability losses to customers. To satisfy the demand for reimbursement for such losses and to bond their services further, partners generally extend their liability to tangible assets held outside the organization (that is, they contract for unlimited liability), or they purchase insurance against liability losses to clients. Such use of

---

[17] Fama & Jensen, *supra* note 2.

unlimited liability and insurance is consistent with the proposition that the dominant asset in a professional partnership is the inalienable human capital of the partners.

2. Unlike professional partnerships, the proprietorships, partnerships, and closed corporations observed in small-scale production activities commonly have mechanisms for transferring residual claims to the cash flows generated by assets other than human capital. Buy-out provisions with internal pricing rules for residual claims and first refusal rights are examples of such mechanisms. Moreover, the residual claims of these organizations are similar in other respects to those of professional partnerships, for example, restriction of the residual claims to important decision agents and periodic renegotiation of salaries to reflect variation through time in the contribution of human capital to net cash flows.

3. Most important, professional partners drop the limited horizon feature of their residual claims when there are substantial assets in the organization in addition to the human capital of existing partners. For example, a departing partner is generally compensated for his share in assets, such as cash and accounts receivable. More interesting, professional partnerships sometimes have devices for compensating a retiring partner for information about his clients that he passes along to remaining partners. Such payments for information reduce the incentives of partners to take actions that substitute near-term cash flows for long-term cash flows in a manner that inhibits organizational survival. It is also interesting that organizations in business and financial consulting that were once professional partnerships with limited horizon residual claims are tending to reorganize as open corporations. We hypothesize that this is largely caused by the pressure to transfer the rights to valuable nonhuman capital assets owned within the organization from one generation of residual claimants to the next.

## B. Financial Mutuals

A common form of organization in financial activities is the mutual. In some financial activities, including life insurance, casualty insurance, and personal savings, mutuals exist side by side with open corporations, and there is no obvious tendency for one form of organization to dominate. Mutuals are dominant among investment mutual funds, but commercial banks are always corporations. Our task is to explain why mutuals survive in some financial activities but not in others.

1. *The Control Function of Redeemable Claims.* An unusual characteristic of mutuals is that the residual claimants are customers, for example, the policyholders of mutual insurance companies, the depositors

of mutual savings banks, and the shareholders of mutual funds. However, the unique characteristic of the residual claims of mutuals, which is important in understanding their survival value, is that the residual claims are redeemable on demand. The policyholder, depositor, or shareholder can, at his initiative, turn in his claim at a price determined by a prespecified rule. For example, the shareholder of an open-end mutual fund can redeem his claim for the market value of his share of the fund's assets, while the whole life or endowment insurance policyholder, like the shareholder of a mutual savings bank, can redeem his claim for its specified value plus accumulated dividends.

There is a special form of diffuse control inherent in the redeemable claims of financial organizations. The withdrawal decisions of redeemable claim holders affect the resources under the control of the organization's managers, and they do so in a more direct fashion than customer decisions in nonfinancial organizations. The decision of the claim holder to withdraw resources is a form of partial takeover or liquidation which deprives management of control over assets. This control right can be exercised independently by each claim holder. It does not require a proxy fight, a tender offer, or any other concerted takeover bid. In contrast, decisions of customers in open nonfinancial corporations, and the repricing of the corporation's securities in the capital market, provide signals about the performance of its decision agents, but without further action, either internal or from the corporate takeover market, the judgments of customers and of the capital market leave the assets owned within the organization under the control of the managers.

2. *The Limitations of Redeemable Claims.* Redeemable claims are not an efficient general financing instrument for nonfinancial organizations. Giving every claim holder the right to force contractions of assets would impose substantial costs on nonfinancial activities. For example, nonfinancial corporations typically have large demands for organization-specific assets that have lower value to other organizations. Substantial costs would be incurred in forced sales of such illiquid assets to accommodate redemptions of claims. In contrast, a financial organization purchases and sells financial assets to meet purchases and redemptions of claims. This is accomplished at low cost because financial assets are not organization specific and can be traded with low transactions costs.

There is a more subtle problem with redeemable residual claims in nonfinancial activities. The pricing rule used to redeem claims preempts development of an outside secondary market for the claims. No one will buy at a price higher than the redemption price or sell at a lower price. The absence of secondary markets for the redeemable claims of financial organizations is no problem since redemption price rules (for example, the net asset value rule for mutual fund shares) can be based on prices of

financial assets quoted in the capital market. In contrast, the residual claims of nonfinancial organizations are claims on uncertain future cash flows. Without a secondary market for the claims, accurate and inexpensive external indexes of their value would not exist, and any internal redemption pricing rule would be costly or arbitrary.

3. *Corporate Financial Organizations.*    Our analysis should also explain why some financial organizations are mutuals and others are open corporations. The theory predicts that more of the business of financial mutuals is management of portfolios of financial assets whereas corporate financial organizations are more involved in business activities requiring organization-specific assets that are expensive to trade and that generate uncertain future net cash flows that are not easily priced.

Observation of different financial organizations is roughly consistent with these hypotheses. Most investment mutual funds manage portfolios of traded securities. The funds are open-end mutuals with redeemable residual claims, except for a handful of closed-end funds organized as open corporations with nonredeemable common stock residual claims. Consistent with our hypothesis, the closed-end funds often hold assets such as real estate or shares in new ventures that are expensive to value and to trade, though this is not universal.[18]

Commercial banks are required by law to be corporations. Our analysis suggests that they would be corporations in the absence of the requirement. A major part of bank business is providing transaction services. Depositors pay for these services directly or by forgoing returns on deposits. The primary assets of commercial banks are short-term loans. Granting and renewing these loans involves monitoring the borrowers and certifying credit worthiness—a service for which the borrowers pay. The capital value of the stochastic net cash flows from services to depositors and borrowers would not easily be captured in the internal pricing rule of a redeemable residual claim.

What survives in commercial banking is a contract structure involving deposits that, like all redeemable claims, allow the depositors to affect the resources under management control. Consistent with our model, variation in deposits is met by purchases and sales of government and private bonds traded at low cost in secondary markets. Since depositors do not have residual claims on net cash flows from service and other activities, redemption of deposits does not require internal valuation of these net cash flows. The rights to the residual net cash flows are assigned to

---

[18] See Rex Thompson, Capital Market Efficiency, Two-Parameter Asset Pricing and the Market for Corporate Control: The Implications of Closed-End Investment Company Discounts and Premiums (1978) (Ph.D. dissertation, Univ. Rochester, Graduate School of Management).

TABLE 1
BUSINESS RECEIPTS AND LONG-TERM NONFINANCIAL ASSETS OF CORPORATE AND MUTUAL
FINANCIAL ORGANIZATIONS, SELECTED YEARS

|  | 1967 | 1969 | 1971 | 1973 | 1975 |
|---|---|---|---|---|---|
| Business receipts as a percentage of total receipts: |  |  |  |  |  |
| Corporate commercial banks | 13.6 | 12.1 | 14.0 | 12.0 | 8.3 |
| Savings and loans | 4.7 | 4.7 | 6.3 | 5.4 | 5.6 |
| Mutual savings banks | 2.9 | 3.0 | 3.1 | 2.8 | 3.1 |
| Corporate life insurance | 82.7 | 82.7 | 83.0 | 82.0 | 81.0 |
| Mutual life insurance | 72.9 | 72.6 | 72.9 | 72.1 | 72.1 |
| Corporate casualty insurance | 91.5 | 89.2 | 89.7 | 87.7 | 87.1 |
| Mutual casualty insurance | 94.0 | 93.0 | 92.7 | 92.0 | 90.1 |
| Long-term nonfinancial assets as a percentage of total assets: |  |  |  |  |  |
| Corporate commercial banks | 2.4 | 2.7 | 3.0 | 3.2 | 3.0 |
| Savings and loans | 2.4 | 2.4 | 2.4 | 2.4 | 2.4 |
| Mutual savings banks | 1.2 | 1.1 | 1.2 | 1.6 | 1.7 |
| Corporate life insurance | 4.9 | 6.1 | 5.4 | 5.4 | 6.5 |
| Mutual life insurance | 2.8 | 3.1 | 3.2 | 3.3 | 3.5 |
| Corporate casualty insurance | 5.3 | 7.6 | 9.0 | 9.5 | 9.5 |
| Mutual casualty insurance | 3.6 | 3.9 | 3.7 | 3.9 | 3.6 |

SOURCE.—U.S. Internal Revenue Service, computer tape of corporate statistics of income. Business receipts are revenues other than interest, dividends, and capital gains. Policy premiums are included in business receipts for insurance companies.

common stock. Since the common stock is not redeemable, there are incentives for development of a secondary market. The residual claims against uncertain future net cash flows are then priced more effectively than would be the case with redeemable residual claims for which there would be no secondary market. Such mixed capital structures, with fixed value redeemable claims (policies or deposits) and nonredeemable common stock residual claims, are also characteristic of the savings banks and insurance companies organized as open corporations.

Our analysis should also explain the differences between the corporate and mutual organizations observed in the same financial activity, for example, life insurance or personal saving. Relative to the mutuals, corporate financial organizations should be more involved in business activities other than management of financial assets, and these business activities should involve relatively more nonfinancial assets that can only be varied with large costs. The data on the business receipts (revenues other than interest, dividends, and capital gains) and long-term nonfinancial assets of banks and life insurance companies in Table 1 are consistent with these

hypotheses. Corporate commercial banks have more business receipts relative to total receipts and more long-term nonfinancial assets relative to total assets than mutual savings banks or savings and loan associations. More interesting, savings and loans, which are sometimes corporations, have relatively more business receipts and long-term nonfinancial assets than mutual savings banks. Likewise, corporate life insurance companies have higher ratios of business receipts to total receipts and higher ratios of long-term nonfinancial assets to total assets than mutual life insurance companies.[19]

The data for casualty insurance organizations are less supportive. Consistent with our analysis, mutual casualty companies show lower ratios of long-term nonfinancial assets to total assets than corporate casualty companies. However, contrary to our analysis, the mutuals have higher ratios of business receipts to total receipts.[20]

Finally, an interesting organizational experiment is taking place in the banking sector. Although commercial banks are required to be corporations, regulations restricting commercial banks and savings banks to different activities are being relaxed. The direction is toward allowing savings banks to provide services such as checking privileges and short-term business loans, previously restricted to commercial banks. If the dominance of the corporate format in commercial banking is not the consequence of regulation, then as savings banks become involved in the service activities of commercial banking, they will tend to organize as corporations. On the other hand, if commercial banking services can be provided at lower prices with the mutual format, corporate commercial banks will not survive when mutual savings banks are allowed to compete with them.

### C.  Nonprofit Organizations

The familiar economic analysis of the entrepreneurial firm is of little help in explaining the dominance of nonprofits in some activities, such as religion, education, research, and classical music, but not in others, including automobile manufacturing, legal services, and popular music. We explain the survival of nonprofits in donor-financed activities as an efficient solution to the special agency problem posed by private donations.

---

[19] Because policy premiums are included as business receipts, business receipts are a larger fraction of total receipts for insurance companies than for banks. Nevertheless, comparison of the business receipts of corporate and mutual insurance companies is relevant.

[20] See David Mayers & Clifford W. Smith, Jr., Contractual Provisions, Organizational Structure, and Conflict Control in Insurance Markets, 54 J. Bus. 407–33 (1981), for additional hypotheses regarding contract structures in the insurance industry.

1. *Nonprofit Organizations and Donations.*     Donations per se do not imply dominance for the nonprofit form. When donations are applied directly to well-defined units of output, a for-profit producer perceives them as a reduction in variable costs or as an increase in demand and increases output accordingly. In fact, we observe unit subsidies both in activities organized on a nonprofit basis, for example, educational scholarships, and in activities organized on a for-profit basis, for example, free tickets to sports events for various groups.

However, some donors wish to provide general donations to particular producers (churches, universities, etc.) rather than unit subsidies. Such unrestricted donations pose agency problems for any organization with residual claimants. Residual claimants contract for rights to net cash flows. When activities are financed in part through donations, part of net cash flow is from resources provided by donors. Contracts that define the share of residual claimants in net cash flows are unlikely to assure donors that their resources are protected from expropriation by residual claimants. One solution to this agency problem is to have no alienable residual claims and to contract with donors to apply all net cash flows to output. Thus, our hypothesis is that the absence of residual claims avoids the donor-residual claimant agency problem and explains the dominance of nonprofits in donor-financed activities.[21]

The absence of alienable residual claims in nonprofits does not mean that residual risk is not borne. When net cash flows are used to expand outputs or to lower the prices of outputs, part of the risk of net cash flows is borne by consumers and part by the factors used to produce the outputs. Thus, residual net cash flows are allocated, but there are no specific residual claimants with alienable property rights in net cash flows. Moreover, the absence of residual claims does not mean that nonprofits make no profits. It means that alienable claims to profits do not exist.

Donations can substitute for the resources provided by residual claimants to purchase assets that are optimally owned rather than rented. When held as endowment, donations also help to bond contracts with

---

[21] Henry B. Hansmann, The Role of Nonprofit Enterprise, 89 Yale L. J. 835 (1980), analyzes the nonprofit organization in detail, but he tends to attribute the nonprofit form more to the nature of products than to the agency problems of donations. He treats donors as customers and looks for product characteristics that would make for "contract failure" in a for-profit framework. For example, charity is delivered to third parties, and the customer (donor) has difficulty verifying delivery. Hansmann also argues that the nonprofit form is attractive for high technology goods (because the customer has difficulty verifying quality) and public goods. However, his approach predicts wider dominance for nonprofits (for example, all high technology or public goods) than is observed. The hypothesis that the nonprofit form is related to donor financing is more promising.

other agents in the organization. From a survival viewpoint the advantage of donations over resources provided by residual claimants is that donors forgo claims on their donations and on the returns earned on the donations, and this tends to allow the organization to deliver its products at lower prices.

Our nonprofit hypothesis deals only with activities financed by donations. Such donor-financed activities are dominated by nonprofits, for example, private universities, churches, hospitals, charities, and cultural performing groups (symphony orchestras, ballet companies, and opera companies). However, the limited scope of the hypothesis means that it cannot explain the nonprofits observed in activities where donations play no role, for example, country clubs.

2. *Other Explanations for Nonprofits.* One criticism of our hypothesis about the causal relation from donations to the nonprofit form is that it ignores the difficulty of measuring and selling the outputs of, for example, churches. The inference is that this explains the nonprofit form in these activities. It is difficult to measure all the things one gets from religion, education, research, or cultural activities. However, the same is true of products such as rock music and legal or psychiatric services marketed by organizations that have residual claims. Moreover, if donations disappeared, for-profit organizations, or more precisely organizations that have alienable residual claims, would arise to supply religion, research, and education. Some for-profit organizations supply these services now. For-profit educational organizations and research groups sell definable parts of their outputs; tuition for education and royalties to patents are examples. For-profit churches might sell ordinations, indulgences, or admission to services. Consistent with our hypothesis, when education and research are provided by organizations that have alienable residual claims, these organizations are not also financed with donations.

Some argue that sale of some products and services (for example, religion) is not acceptable and that this explains the nonprofit form in these activities. This is consistent with our hypothesis. When giving outputs away generates more resources through donations than sale, survival dictates the nonprofit form. Thus, universities generally make research freely available because this generates more resources through research grants and other donations than direct sale of the research. Churches usually do not insist on payment of admission charges or member taxes because they attract more total resources through voluntary contributions.

Coldly economic statements like these lead to the criticism that our analysis leaves no room for altruism. The opposite is true. Altruistic internal agents increase the willingness of altruistic customers and donors

to provide resources. In our terms, the altruism of internal agents allows low cost control of agency problems and acts to bond donors and customers against expropriation. Strong tastes for an organization's outputs on the part of internal agents and customers—what we call altruism in the case of nonprofits—contribute to the survival of any organization. All organizations try to develop such brand loyalty, but the nonprofits are especially successful, perhaps because of the nature of their products.

Some readers claim that donors, customers, and internal agents have tastes for the nonprofit form itself in some activities. To explain the complete dominance of nonprofits in an activity, however, this approach requires uniformity of tastes. If subgroups of customers, internal agents, and donors have no preference for the nonprofit form, we would expect more competition among profit and nonprofit organizations in donor-financed activities.

Finally, tax concessions are important to some nonprofits. However, the major activities dominated by nonprofits, such as religion, private education, research, hospital care, and certain cultural activities, were dominated by nonprofits before taxes were a major issue.[22] Our hypothesis about the relation between unrestricted donations and the nonprofit form provides a more consistent explanation of the historical dominance of nonprofits in these activities. On the other hand, tax exemptions probably explain the nonprofits in activities where private donations are not a factor, including nursing homes, homes for the elderly, and private nursery schools.

3. *The General Control Problem in Nonprofits.* The donors of nonprofits have agency problems with internal decision agents similar to those faced by residual claimants in other organizations, such as open corporations and financial mutuals, where important decision managers do not bear a major share of the wealth effects of their decisions. We argue in "Separation of Ownership and Control"[23] that, like all other organizations characterized by separation of decision management from residual risk bearing, a nonprofit is on stronger footing in the competition for survival when it has a decision system that separates the management (initiation and implementation) and control (ratification and monitoring) of important decisions. For nonprofits the survival value of such decision systems is due to the assurances they provide that donations are used effectively and are not easily expropriated.

For example, like open corporations and financial mutuals, donor

[22] See *id*.
[23] Fama & Jensen, in this issue.

nonprofits have boards of directors (or trustees) with the power to ratify and monitor important decisions and to hire, fire, and set the compensation of important decision agents. The similarities of the decision control systems of nonprofits, financial mutuals, and open corporations, along with the differences due to special agency problems and special features of residual claims (including the absence thereof), are discussed in "Separation of Ownership and Control."

## V.  SUMMARY AND CONCLUSIONS

Most goods and services can be produced by any form of organization. Organizations compete for survival, and the form of organization that survives in an activity is the one that delivers the product demanded by customers at the lowest price while covering costs.

The characteristics of residual claims are important both in distinguishing organizations from one another and in explaining the survival of specific organizational forms in specific activities. We explain the survival of organizational forms largely in terms of the comparative advantages of characteristics of residual claims in controlling the agency problems of an activity. The analysis identifies the underlying characteristics of activities that determine the organizational forms that survive.

### A.  Open Corporations

The common stock residual claims of open corporations are unrestricted in the sense that (1) they are freely alienable, (2) they are rights in net cash flows for the life of the organization, and (3) stockholders are not required to have any other role in the organization. Other things equal, the open corporation is more likely to survive in an activity the greater
1. the benefits of unrestricted risk sharing,
2. the benefits of specialized management,
3. the amount of organization-specific assets to be purchased,
4. the wealth required to bond contractual payoffs, and
5. the lower the cost of separating decision management (initiation and implementation) from decision control (ratification and monitoring).

For example, these factors favor the open corporate form when the technology in an activity implies economies of scale that involve (*a*) large aggregate residual risks to be shared among residual claimants, (*b*) large demands for specialized decision agents throughout the organization, and (*c*) large demands for wealth from residual claimants to bond contracts and to purchase organization-specific assets. Economies of scale are also likely to imply organizations that are complex in the sense that valuable specific knowledge—knowledge that is expensive to transfer across

agents—is widely diffused among agents.[24] Such complexity tends to favor unrestricted common stock residual claims which allow specialization of management and delegation of decision functions to agents with valuable relevant knowledge.

The benefits of unrestricted common stock residual claims in activities where optimal organizations are large and complex offset the agency costs resulting from the separation of decision functions and residual risk bearing. In "Separation of Ownership and Control" we contend that these agency costs are controlled by decision structures that separate the management and control of important decisions.

### B.  Proprietorships, Partnerships, and Closed Corporations

In a fictional world where contracts with decision agents were cost-lessly written and enforced, separation and specialization of decision and risk-bearing functions would involve no agency costs, and most if not all organizations would have unrestricted residual claims. However, actual organizations can realize the benefits of unrestricted residual claims only by incurring costs to control agency problems between specialized decision agents and specialized residual risk bearers. As a consequence, it is advantageous in some activities to trade the benefits of unrestricted common stock residual claims for the low-cost control of agency problems in the decision process obtained when residual claims are restricted to important decision agents. This restriction is a common characteristic of the residual claims of proprietorships, partnerships, and closed corporations. Other things equal, these organizations with their restricted residual claims are more likely to survive in activities where the costs of separating decision management from decision control are high. They are also more likely to survive when there are no important economies of scale and thus (a) no large demands for unrestricted risk sharing and specialized decision skills, and (b) no large demands for wealth from residual claimants to bond contracts and purchase organization-specific assets.

### C.  Special Forms of Residual Claims

Organizations such as professional partnerships, financial mutuals, and nonprofits have residual claims with unique characteristics that we explain as devices for controlling special agency problems.

1.  *Professional Partnerships.*   These are characterized by (1) restriction of residual claims to major decision agents, (2) periodic renegotiation

---

[24] The role of specific knowledge is discussed in Fama & Jensen, Separation of Ownership and Control, in this issue.

of partner shares in net cash flows (flexible sharing rules), and (3) inalienable residual claims in net cash flows with horizons that are often limited to a partner's period of service in the organization. Professional partnerships are more likely to survive in an activity when

1. valuable specific knowledge relevant to both the management and control of decisions is combined and diffused among agents,

2. there are no strong demands for organization-specific tangible assets, and

3. the benefits from consulting and mutual monitoring among decision agents are high.

These characteristics are observed in professional service activities (law, public accounting, and business consulting) where (1) restricting residual claims to important decision agents helps control the agency problems caused by delegating combined decision management and control rights with respect to cases, audits, and so forth, to agents with relevant specific knowledge; (2) the primary asset of the activity is professional human capital; and (3) mutual monitoring and consulting among agents are important to maintain the value of human capital, which is sensitive to performance.

2. *Financial Mutuals*.  The distinguishing characteristic of the residual claims of financial mutuals is that the policyholder, depositor, or shareholder can sell his claim to the organization on demand at a price determined by a rule. The decision to withdraw resources by the holder of a redeemable claim is a form of partial takeover or liquidation that deprives management of control over assets. This mechanism for decision control can be exercised independently by each claim holder. It does not require a proxy fight, a tender offer, or any other concerted takeover bid. Mutuals are more likely to survive in an activity the lower the cost

1. of expanding and contracting assets and

2. of obtaining accurate indices of asset values.

These conditions occur in financial organizations where assets are primarily the securities of other organizations. Redeemable residual claims are a low-cost mechanism for controlling agency problems between the residual claimants and the decision agents of financial mutuals because accurate and inexpensive indexes for asset values are available and the assets are traded with low transactions costs. Redeemable claims are a high-cost mechanism for decision control in activities that involve large amounts of assets not traded in secondary markets. Redeemable residual claims are also inefficient in activities that involve large amounts of lumpy or organization-specific assets that can be varied only with large costs.

3. *Nonprofits*.  The nonprofit organization is characterized by the

absence of alienable residual claims to net cash flows and contractual constraints on the distribution of net cash flows. Inalienable residual claims are vested in a board of trustees and net cash flows are committed to current and future output. Nonprofits are more likely to survive in an activity

1. the greater is the potential supply of donations and

2. the lower is the cost of separating decision management from decision control.

The nonprofit organization is a solution to the agency problem posed by donations. When the activities of an organization are financed in part through donations, part of stochastic net cash flow is due to the resources provided by donors. Contracts that define the share of residual claimants in net cash flows are unlikely to assure donors that their resources are protected against expropriation by residual claimants. One solution to this agency problem between donors and residual claimants is to have no residual claimants and to contract with donors to apply net cash flows to future output. The absence of alienable residual claims means that decision managers in nonprofits do not bear the wealth effects of their decisions. As in other organizations where residual risk bearing and decision management functions are separated, the resulting agency problems in the decision process are controlled by decision structures that separate the management and control of important decisions.

## APPENDIX

### BIBLIOGRAPHY

Arrow, Kenneth J. "The Role of Securities in the Optimal Allocation of Risk Bearing." *Review of Economic Studies* 31, no. 86 (1964): 91–96.

Berle, Adolf A., and Means, Gardiner C. *The Modern Corporation and Private Property*. New York: Macmillan Publishing Co., 1932.

DeAngelo, Linda E. "Auditor Size and Audit Quality." *Journal of Accounting and Economics* 3, no. 3 (December 1981): 183–200.

Debreau, Gerard. *Theory of Value*. New York: John Wiley & Sons, 1959.

Fama, Eugene F. *Foundations of Finance*. New York: Basic Books, 1976.

Fama, Eugene F. "The Effects of a Firm's Investment and Financing Decisions on the Welfare of Its Security Holders." *American Economic Review* 68, no. 3 (June 1978): 272–84.

Fama, Eugene F., and Jensen, Michael C. "Separation of Ownership and Control." *Journal of Law & Economics* 26 (June 1983): 301–25.

Fama, Eugene F., and Jensen, Michael C. "Organizational Forms and Investment Decisions." Managerial Economics Research Center Working Paper no. MERC 83-03. Rochester, N.Y.: University of Rochester, Graduate School of Management, 1983.

Furubotn, E. G., and Pejovich, S. "Property Rights, Economic Decentralization

and the Evolution of the Yugoslav Firm, 1965–1972." *Journal of Law & Economics* 16 (October 1973): 275–302.

Hansmann, Henry B. "The Role of Nonprofit Enterprise." *Yale Law Journal* 89, no. 5 (April 1980): 835–901.

Jensen, Michael C. "Organization Theory and Methodology." *Accounting Review* 50 (April 1983).

Jensen, Michael C., and Meckling, William H. "Theory of the Firm: Managerial Behavior, Agency Costs and Ownership Structure." *Journal of Financial Economics* 3, no. 4 (October 1976): 305–60.

Jensen, Michael C., and Meckling, William H. "Rights and Production Functions: An Application to Labor-managed Firms and Codetermination." *Journal of Business* 52, no. 4 (October 1979): 469–506.

Klein, Benjamin; Crawford, Robert; and Alchian, Armen A. "Vertical Integration, Appropriable Rents, and the Competitive Contracting Process." *Journal of Law & Economics* 21, no. 2 (October 1978): 297–326.

Mayers, David, and Smith, Clifford D., Jr. "Contractual Provisions, Organizational Structure, and Conflict Control in Insurance Markets." *Journal of Business* 54, no. 3 (July 1981): 407–33.

Mayers, David, and Smith, Clifford W., Jr. "On the Corporate Demand for Insurance." *Journal of Business* 55, no. 2 (April 1982): 281–96.

Reagan, Patricia B., and Stulz, Rene M. "Risk Bearing, Labor Contracts and Capital Markets." Managerial Economics Research Center Working Paper Series no. MERC 82-19. Rochester, N.Y.: University of Rochester, Graduate School of Management, November 1982.

Smith, Adam. *The Wealth of Nations,* 1776. Edited by Edwin Cannan, 1904. Reprint. New York: Modern Library, 1937.

Thompson, Rex. "Capital Market Efficiency, Two-Parameter Asset Pricing and the Market for Corporate Control: The Implications of Closed-End Investment Company Discounts and Premiums." Ph.D. dissertation, University of Rochester, Graduate School of Management, 1978.

# AGENCY PROBLEMS, AUDITING, AND THE THEORY OF THE FIRM: SOME EVIDENCE*

ROSS L. WATTS and JEROLD L. ZIMMERMAN
University of Rochester

## I. INTRODUCTION

Recent developments in the theory of the firm emphasize the impor-
tance of monitoring the performance of parties to the firm.[1] Jensen and
Meckling[2] hypothesize that an audit is one type of monitoring activity that
increases the value of the firm. An audit by someone independent of the
manager reduces the incentive problems that arise when the firm manager
does not own all the residual claims on the firm. If Jensen and Meckling's
hypothesis is correct, independent audits are expected in the earliest firms
where the manager did not supply all the capital.

On the other hand, a leading auditing text suggests that the appearance
of independent audits is more recent and is the product of government
fiat: "One of the earliest steps in recognition of the need for audits oc-
curred in England with the passage of the Registered Companies Act of
1862. The Act required that the financial statements of joint stock com-
panies be audited by a person independent of management, and thereby
greatly enhanced the status of professional auditors as well as the growth
of that profession."[3] This opinion, that independent audits arose because

---

* Associate Professors. The authors gratefully acknowledge the comments on earlier
versions of this manuscript by Eugene Fama, Nicholas Gonedes, Robert Holthausen,
Michael Jensen, Richard Leftwich, Clifford Smith, Lee Wakeman, and Jerold Warner.

[1] Armen A. Alchian & Harold Demsetz, Production, Information Costs, and Economic
Organization, 62 Amer. Econ. Rev. 777 (1972); Michael C. Jensen & William H. Meckling,
Theory of the Firm: Managerial Behavior, Agency Costs and Ownership Structure, 3 J.
Financial Econ. 305 (1976); Eugene F. Fama, Agency Problems and the Theory of the Firm,
88 J. Pol. Econ. 288 (1980); Eugene F. Fama & Michael C. Jensen, Separation of Ownership
and Control, 26 J. Law & Econ. 301 (1983); Eugene F. Fama & Michael C. Jensen, Agency
Problems and Residual Claims, 26 J. Law & Econ. 327 (1983).

[2] Jensen & Meckling, *supra* note 1, at 338–39.

[3] Howard F. Stettler, Auditing Principles 20–21 (4th ed. 1977). Note that the statement

they were specifically required by government regulation, is echoed in a recent congressional staff report that claims that the U.S. securities acts "created a need for . . . independent auditors."[4]

To provide evidence on the importance of the independent audit to the firm, this paper investigates the history of the business corporation. The organization and use of audits by English and U.S. business corporations are examined, from the days of English merchant guilds to the time audits were required by law. We find that the audit existed early in the development of business corporations (1200) and evolved gradually into the type of audit required by the first English companies act (1844). Further, it appears that audits were widespread among early business corporations. This evidence suggests that monitoring of performance is important, if not crucial, to the formation of firms. The long survival of auditing suggests it is a part of the efficient technology for organizing firms.[5]

The audits of the early corporations were conducted by directors or shareholders. The use of outside professional auditors did not become common until the latter half of the nineteenth century in the United Kingdom and the early part of this century in the United States. However, when the use of professionals became common, corporate audits were not generally required by law in either country. This suggests the use of professional auditors was due to changes in the market for auditing.

The papers that stress the importance of monitoring in the theory of the firm analyze the incentive (agency or control) problems that exist in different types of firms (partnerships, nonprofit firms, corporations, etc.). Firms are sets of contracts among the factors of production,[6] and different sets of contractual arrangements (for example, alternative property right structures) provide different incentives for opportunistic behavior by the contracting parties. This opportunistic behavior reduces the total product of the firm and hence its value. Jensen and Meckling point out that in markets characterized by rational expectations the cost of the opportunistic behavior is borne by the offending party. This provides incentives for that party to write contracts which restrict his opportunistic behavior.

Corporate managers or promoters write contracts (by-laws, charters, indenture agreements, compensation plans) that reduce their opportun-

that the 1862 Companies Act required independent auditing is false. That act neither required audits nor specified that the auditors be independent. That act contained an *optional* model set of articles that provided for audits but did not specify they be by outsiders.

[4] Staff of Senate Subcomm. on Reports, Accounting, and Management of the Committee on Government Operations, 94th Cong., 2d Sess., The Accounting Establishment: A Staff Study 1 (Comm. Print 1976).

[5] Armen A. Alchian, Uncertainty, Evolution, and Economic Theory, 58 J. Pol. Econ. 211 (1950).

[6] Jensen & Meckling, *supra* note 1.

ism. Enforcement of the contract requires monitoring of management's activities and it is hypothesized by Jensen and Meckling that this is a role of auditing. Another enforcement mechanism discussed by Jensen and Meckling is the posting of bonds by managers. The evidence reviewed below reveals that, like auditing, bonding has been used throughout the history of the business corporation.

An audit will be successful in changing expectations and hence reducing the opportunistic behavior costs (agency costs) borne by the manager only if it is expected that the auditor will report some discovered breaches of contract. The probability that the auditors will report a discovered breach is effectively the auditing profession's definition of independence.[7] So the preceding condition can be rephrased to state that an audit will reduce the agency costs borne by the manager only if the market expects the auditor to have a nonzero level of independence. The incentives of the early business corporations' auditors to be independent are reviewed in this paper. Further, we analyze the changes in the mechanics of the audit and in the auditor's incentive to be independent as the business corporation evolved.

The evidence analyzed in this paper is qualitative in nature and does not lend itself to hypothesis testing using large samples and econometric methods. Random samples of company charters, by-laws, or minutes cannot be drawn since what exist are chiefly the records of successful companies. Although audits existed in every firm whose original documents we examined (the English Merchant Adventurers, regulated companies, and joint stock companies from the fifteenth to the seventeenth centuries), there may have been other companies that did not have audits and whose records were not preserved. In spite of this ex post selection bias, we believe that this historical evidence provides interesting insights on the theory of the firm and on the origins of the independent audit.

Section II presents evidence that audits were common in early English business corporations, the medieval craft and merchant guilds, and regulated and early joint stock companies. It presents evidence also that the auditors had incentives to be independent. This evidence is consistent with the hypothesis that efficient, agency-cost-minimizing contractual arrangements (for example, the organizational forms) include monitoring and is inconsistent with the explanation that independent corporate auditing is the consequence of companies or securities acts. Section II describes also how the audit evolved in response to changes in organizational forms. Early audits were performed by committees of owners

---

[7] Ross L. Watts & Jerold L. Zimmerman, Auditor Independence and Scope of Services (July 1982) (unpublished manuscript, Univ. Rochester, Grad. School Management).

(members). The composition and size of these committees changed as the organization's size and structure changed.

Section III describes the development of the independent professional audit firm and its substitution for the internal shareholder audit committee. In addition the section explores explanations for that substitution. Section IV summarizes the conclusions.

## II. Auditing in the Early Corporation

In this section we provide a description of each type of business corporation to demonstrate that incentive problems existed in those corporations. Then we give the evidence that auditing and bonding existed for the type of corporation along with the evidence that the audits were designed to reduce incentive problems. Finally, we discuss the incentives of the auditors to be independent.

### A. English Merchant Guilds

1. *Description.* Merchant guilds appeared in England shortly after the Norman Conquest (A.D. 1066).[8] The guild arose to protect the prosperity of the merchants after the conquest by forming cartels to monopolize trade.[9] They did not exist before the conquest because there was very limited trade.[10] A guild was chartered by the crown and given a monopoly to trade within its own particular town. Each guild member traded on his own account or in partnership with other guild members. Members of a guild were not allowed to enter into partnership with nonmembers. Further, a guildman had to share any purchase with other guild members who wished to participate in the venture, at the same price.[11]

The guilds are among the earliest examples of incorporation. Most guilds held property and incorporation reduced the transactions costs of

---

[8] 1 Charles Gross, The Gild Merchant 4 (1890); 1 W. Ashley, An Introduction to English History and Theory 71 (1923); 1 William Robert Scott, The Constitution and Finance of English, Scottish and Irish Joint-Stock Companies to 1720, at 7 (1912).

[9] Robert E. Ekelund & Robert D. Tollison, Mercantilist Origins of the Corporation, 11 Bell J. Econ. 715 (1980). Besides functioning as a cartel, guilds performed other functions including social, religious, and municipal government. Although these other functions for some guilds might be more important than the cartel-enforcing, mercantile functions, it is these latter functions that are the focus of the subsequent analysis. Carlo M. Cipolla, Before the Industrial Revolution: European Society and Economy, 1000–1700 (2d ed. 1980).

[10] Gross, *supra* note 8, at 3–4.

[11] In the twelfth to fourteenth centuries, the merchant guilds in the larger towns started to specialize by craft or ware. The organization of these craft guilds (sometimes called mercantile societies or companies of merchants) is similar to that of the merchant guilds. Gross, *supra* note 8, at 106–57; and Scott, *supra* note 8, at 8.

transferring the property.[12] If the guild itself was not able to hold property there were substantial costs in changing titles as individual members died and new members were admitted. Hence the idea arose that the organization itself was "to continue indefinitely as the owner of the premises devised."[13] The idea that the whole body could act and that the act was separate from that of a member was expressed in the development of the common seal to be used as evidence that the corporation itself was acting.[14] Guilds were expressly incorporated as early as the reign of Richard II (1367–1400).[15]

The members of the guild provided resources administered by the officers of the guild.[16] Contractual arrangements (for example, charters and by-laws) existed to define and restrict the actions of the contracting parties, including the officers. The major officers generally were an alderman or master and his associates, called stewards, skevins, or wardens. The number of stewards generally varied from two to four and in some places the stewards fulfilled the alderman's role.[17] The alderman and stewards administered the revenues from entrance fees, assessments, and tolls. Many guilds had lands and tenements.[18] In addition, the officers often traded on behalf of the members and conflicts of interests between manager and owners were present. Given that, it is not surprising the manager's acts were restricted by the contracts. The alderman and stewards had specific duties and roles to fulfill.[19] Further, given the restrictions on the guild officers' actions, we would expect to observe that the officers' actions would be monitored.

2. *The Existence of Auditing and Bonding.*   The constitution of the merchant guild at Ipswich in 1200 requires that the alderman "on oath shall make due return, annually before the bailiffs and coroners (of the town) of all profits arising during the year" from the guild's monopoly trading of stone and marble.[20] The merchant guild at Bury St. Edmunds had, by 1304, provision for an annual audit.[21]

---

[12] Gross, *supra* note 8, at 99.

[13] Scott, *supra* note 8, at 3.

[14] *Id.*

[15] Gross, *supra* note 8, at 99.

[16] Although it is difficult to assess the typical proportion of a member's wealth administered by the guild officials, some guilds appear to have owned substantial property. Scott, *supra* note 8, at 6; and Gross, *supra* note 8, at 14, 151, 159–60.

[17] Gross, *supra* note 8, at 26–27.

[18] *Id.* at 28.

[19] *Id.* at 23–35.

[20] *Id.* at 25.

[21] This reference is to the original medieval Latin text of documents supplied by 2 Gross,

Several craft guilds and companies of merchants were audited annually by committees of members.[22] The records of the Worshipful Company of Grocers of the City of London and the Worshipful Company of Pewterers of the City of London indicate, in 1346 and 1546, respectively, that the accounts of those two companies were audited annually. Similarly, the accounts of the retiring Warden of the Worshipful Company of Carpenters were audited each year from the fifteenth century on by a committee of past wardens and other members and based on their report, the balance, if any, of his bond[23] returned to the retiring Warden.[24] The evidence suggests that the audits were not superficial and were not merely counts of cash or assets on hand. Expenditures were examined in detail. Boyd cites several examples of the auditors refusing to certify and disallowing various charges.

The audit of the guild appears designed to monitor the managers' contracts. It came at the end of the managers' tenure and was designed to check for unauthorized expenditures. Further, it also appears to have been designed to check for other breaches of contracts. In at least one case, the auditors of a craft guild explicitly fined the master and wardens for breaching certain ordinances (that is, breaching their contract).[25] The prevalence of guild audits[26] and bonding arrangements is consistent with the existence of conflicts of interest creating incentives for bonding and monitoring technologies to be devised.

Guild audit committees usually consisted of four guild members and occasionally public officials.[27] Although these auditors did not specialize in auditing as a full time career, they still had the responsibility to monitor certain managerial functions and to check for breaches of contract.

3. *The Auditors' Incentives to Be Independent.*   An auditor could be heavily fined for not completing the audit in due time.[28] In addition to that

---

*supra* note 8, at 34, to support his statements in vol. 1. We obtained the details of the audit by having the Latin translated.

[22] Edward Boyd, History of Auditing, in History of Accounting and Accountants 78–88 (Richard Brown ed. 1905).

[23] Bonding is hypothesized by Jensen & Meckling, *supra* note 1, at 325, to be a way to reduce the costs that arise from the conflicting of interests.

[24] Boyd, *supra* note 22, at 86–88.

[25] *Id.* at 81.

[26] The development of the incorporated borough is closely related to the development of the merchant guild; the officers of the borough tend to be drawn from guild members (see 1 Gross, *supra* note 8). Observing borough audits similar to those of the guilds (see Boyd, *supra* note 22) suggests the type of audit described above was a widespread phenomenon in the fourteenth, fifteenth, and sixteenth centuries.

[27] Initially, the auditors were compensated in goods (for example, ale) and later were paid. In the Spanish Company (a regulated company circa 1600) the eight auditors were paid twenty shillings. Pauline Croft, The Spanish Company 79–80 (1973).

[28] Boyd, *supra* note 22, at 81.

cost, lack of performance and independence affected the auditor's reputation and, in the extreme, caused loss of his guild membership and his share of the guild's monopoly profits. In one case, at least, the guild auditors had to own property.[29] Presumably this requirement made it easier to recover damages against them if they breached the contract, thereby providing the auditors with further incentives to report a breach, if one occurred.

The use of an audit committee rather than a single auditor also encouraged auditor performance and independence by making collusion of the manager and auditor more difficult. It is more costly to bribe an entire committee than a single individual.

## B. Merchant Adventurers/Regulated Companies

1. *Description.* Prior to the latter half of the thirteenth century the English export trade was generally conducted by the German Hanse Merchants.[30] However, from that time until the seventeenth century the export of English raw products was conducted through English and foreign merchants who were given a monopoly over the export of raw products (principally wool). Certain towns were specified for the shipping and sales of those products, companies of merchants were eventually formed, and the system facilitated the collection of the royal customs.[31] Each member of the company provided his own capital (inventories and ships) and traded on personal account (or in partnership).

In the fourteenth century, under Edward III's auspices, the cloth industry grew in England and manufactured goods, primarily cloth, began to be exported. In the late fourteenth and early fifteenth century, English merchants involved in this trade applied for charters for organizations to monopolize the trade. In 1391, 1407, and 1408 charters were granted to English merchants trading with Prussia, the Netherlands, and Scandinavia, respectively. The latter two companies became known as the Merchant Adventurers of England and the Eastland Company.[32] In the sixteenth century, companies of merchants trading with other countries were incorporated, including the Spanish Company (1577), the Turkey Company (1581), the Venetian Company (1583), and the Levant Company, a union of the Turkish and Venetian Companies (1592).[33]

---

[29] Gross, *supra* note 8, at 34.

[30] *Id.* at 140.

[31] *Id.*, at 144.

[32] Scott, *supra* note 8, at 8–9.

[33] In exchange for granting and enforcing the cartel, the Crown was often made a partner or given low interest "loan," Ekelund & Tollison, *supra* note 9. The regulated companies passed by-laws governing commerce with the foreign countries. Their rules were specific

There is evidence that some of the members of the early regulated companies were also members of companies of merchants. For example, the minutes of the Merchant Adventurers were kept in the same book as those of the Mercers' Company of London and the Mercers' Hall was Adventurers' headquarters until 1666.[34] These facts suggest some of the merchants who began the Merchant Adventurers came from the Mercers' Company. Given this relationship, it is not surprising that the form of organization of the regulated company is similar to that of the craft and merchant guilds.

By the fifteenth century some guilds had developed a system of administration consisting of a governor and a number of associates, usually a multiple of twelve. For example, the alderman and four associates of the Ipswich guild were replaced by a governor and twenty-four associates.[35] This change is reflected also in the organization of the Merchant Adventurers of England which in its 1505 charter called for the election of a governor and twenty-four associates.[36]

2. *The Existence of Auditing and Bonding.*   Most of the regulated companies appear to have been audited.[37] Before 1632 the Levant Company elected auditors as required (for example, when accounts were presented). After 1632 auditors were elected annually.[38] The Register Book of the Spanish Company includes an entry[39] that at the meeting of the general court (shareholders' meeting) on September 7, 1604, the treasurer "brought in his account" and eight auditors were appointed to examine it.[40] Six were assistants (directors) and two were ordinary members. The Eastland Company also had auditors. A letter from the Company at London to the Company at York in 1674 refers to the auditors' being dissatisfied because the details were not supplied for certain charges.[41]

---

and governed details of family and social life, Scott, *supra* note 8, at 10. Breach of the rules could lead to a member's losing his right to trade with the foreign country.

[34]  Gross, *supra* note 8, at 149.

[35]  Scott, *supra* note 8, at 7.

[36]  *Id.* at 9.

[37]  References to six English regulated companies, two of which later merged, were found. Of the remaining five companies, four were definitely audited. Historical records of the company trading with Prussia were sketchy and contained no mention of an audit.

[38]  Mortimer Epstein, The Early History of the Levant Company 68 (1908).

[39]  Croft, *supra* note 27, at 11 and 79.

[40]  Apparently, accounts were also presented annually to the company before the war with Spain. Croft, *supra* note 27, at xxix, notes that neither the last treasurer before the war, George Hanger, nor his immediate predecessor, Sir John Watts, bothered to present accounts (because the company broke up due to the war). The accounts presented on September 7, 1604, were Hanger's apparently before the war, because the register book also records a resolution warning Watts to present his account with all convenient speed.

[41]  Maud Sellers, The Acts and Ordinances of the Eastland Company 97 (1906).

The codification of the ordinances and decrees of the Merchant Adventurers of Bristol in 1618 includes a requirement that the treasurer present a statement of receipts and payments at the annual meeting and "yield vpp his accompt before Auditors assigned for that purpose."[42]

As in the guilds, the monitoring mechanism of auditing was supplemented by the bonding of officials. The treasurers of the Merchant Adventurers of Bristol,[43] the Spanish Company,[44] the Levant Company,[45] and the Eastland Company[46] were required to post bonds. Further, these bonds were substantial (the Levant Company's was four hundred pounds, the Spanish Company's five hundred pounds, and the Eastland Company's one thousand marks sterling; the Merchant Adventurers do not specify an amount). In addition, officials in towns other than London and in foreign towns were required to post bonds.[47]

The form of the audit of the regulated company is similar to that of the merchant and craft guilds and companies of merchants. A committee of members was elected to audit the treasurer's accounts. One difference was the increase in the size of the committee of auditors, which was probably due to the larger scale of the regulated companies. Not only the general accounts were audited, but also the accounts of the company in other towns.

3. *The Auditor's Incentives to Be Independent.* As in merchant and craft guilds, the auditors could be fined for refusing to act,[48] and the committee form made collusion with the manager more difficult.

It appears from the existence and the similarity of the requirements across all the companies whose records we can obtain that auditing of regulated companies was the usual if not the universal practice. Further, the similarities of audit arrangements (for example, the use of bonding of the officers and the audit to determine how much of the bond to return) and the commonality of membership with the guilds suggests the audit was adopted from the guilds. This in turn suggests that the same mechanisms provided the auditor with incentives to be independent (for example, the use of committees and penalties—including loss of reputation).

---

[42] John Latimer, The History of the Society of Merchant Venturers of the City of Bristol 69 (1903).

[43] *Id.*

[44] Croft, *supra* note 27, at 79.

[45] Epstein, *supra* note 38, at 73.

[46] Sellers, *supra* note 41, at 24.

[47] *Id.* at 25.

[48] *Id.* at xxiii.

### C.   Joint Stock Companies

1. *Description.*   In 1553, one and a half centuries after the appearance of the first regulated company, another form of corporate firm appeared in England, the joint stock company. The first joint stock companies, the Russia Company and the African Adventurers, were formed for overseas trade just as the regulated companies were. However, a major difference between the two was the method of financing the trade. In the regulated company each member supplied his own capital and traded on his own account or in partnership, using his own ships.[49] In the joint stock company the officers of the company traded on behalf of all the members or shareholders. Initially, capital was raised to finance each separate voyage and the proceeds were distributed after the voyage was completed. Several hypotheses can explain this changed method of financing, which occurred while regulated companies continued to be formed after the emergence of the joint stock company (for example, the Turkey Company in 1581, the Venice Company in 1583, and the Levant Company in 1592).[50]

One hypothesis is that the scale of the joint stock company voyages, per se, explains the change. Either the voyages required too large an investment for a single merchant (or a partnership of a few merchants) or the scale of the voyages of regulated companies was too small to allow separate specializations in management and risk bearing.[51] However, it is unlikely that this hypothesis alone can explain the emergence of the joint stock company. The first two joint stock companies were the Russia Company and the African Adventurers, both of which had their first voyages in 1553.[52] It is unlikely that the value of the cargoes of each voyage of the two firms exceeded 5,000 pounds.[53] At the same time, the regulated company, the Merchant Adventurers, shipped 60,000 pounds a year in woolen goods.[54] It would seem a partnership of Merchant Adventurers could have financed a 5,000-pound voyage.

Another hypothesis is that the magnitudes of the voyages combined

---

[49] Thomas S. Willan, The Early History of the Russia Company 1553–1603, at 19 (1956).

[50] Epstein, *supra* note 38.

[51] George J. Stigler, The Division of Labor is Limited by the Extent of the Market, 59 J. Pol. Econ. 185 (1951), discusses the effect of the extent of the market on the degree of specialization.

[52] Scott, *supra* note 8, at 17–22. The original names of the companies were (for the Russia company) "the Merchants Adventurers of England for the discovery of lands, territories, isles, dominions and seignories unknown, and not before that late adventure of enterprise commonly frequented" and (for the African Adventurers) "the Adventurers to Guinie."

[53] Willan, *supra* note 49, at 22.

[54] *Id.*

with the greater risk of the voyages undertaken by the joint stock companies led to risk sharing of the claims on the voyages' proceeds.[55] The
regulated companies evolved to control trade (often already established)
with relatively familiar countries. The routes, the perils faced, the nature
of the trade, and so on, were known. On the other hand, the first voyage
of the Russia Company sought "to discover a northeast passage to the
Indies."[56] The route and the nature of any trade that might result were
unknown. Similarly, the African Adventurers sought to establish trade
with Guinea—a risky undertaking. The Portuguese were certain to attempt to stop such trade and the traders faced danger from the Spanish.[57]
Portfolio theory suggests the efficient allocation of risk bearing leads to
investors being well diversified. Thus, we might expect the higher risk of
the voyages of the Russia Company and the African Adventurers would be
spread across many investors.[58]

The majority of members of the Russia Company were London merchants already engaged in foreign trade (staplers and Merchant Adventurers).[59] Hence, it is not surprising that the 1555 Charter of the Russia
Company was hard to distinguish from that of a regulated company.[60] The
joint stock character of the company was not explicitly recognized in the
charter.[61] Like a regulated company or a guild of the time, the Russia
Company had a governor and twenty-four assistants.

The African Adventurers operated without a charter. It comprised five
chief adventurers who in turn had partners "under" them.[62] The company
made calls on its shares to finance each voyage and then distributed the
proceeds in accordance with the shares. The company probably never
sought a charter because it wished to keep its operations secret to avoid
difficulties with the Portuguese.[63] The company ceased operations in

[55] Fama, *supra* note 1.

[56] Willan, *supra* note 49, at 1.

[57] Scott, *supra* note 8, at 21.

[58] Fama, *supra* note 1.

[59] See Willan, *supra* note 49, at 21. This evidence is inconsistent with the hypothesis
advanced by Ekelund & Tollison, *supra* note 9, who argue that joint stock companies arose
as the efficient method whereby owners of the cartel could sell their rights to the highest
valued users—the most efficient managers. While this is certainly possible, it suggests that
the regulated companies, involving the same merchants as the joint stock companies, would
adopt the more efficient joint stock organization. However, both organizational forms are
observed at the same time, in the same country, and involving the same principals.

[60] Willan, *supra* note 49, at 22.

[61] Scott, *supra* note 8, at 19.

[62] *Id.* at 21, 30.

[63] *Id.* at 21.

1566,[64] and it was not until 1588 that another company was formed to trade with Africa.[65]

2. *The Existence of Auditing and Bonding.* The first joint stock companies were audited. Sources external to the company indicate that as early as the 1580s, and afterwards, the Russia Company was audited annually.[66] The Court of Minutes of the East India Company for the early years of its existence (September 1599–August 1605) includes resolutions for the appointment of auditors every time accounts were presented to the general court (that is, the shareholders' meeting). The first such resolution is dated December 31, 1600, and all four auditors appointed were directors.[67]

The accounts of the early joint stock companies were audited by a committee of shareholders (members) and/or directors. This practice continued into the eighteenth and nineteenth centuries. Some companies had provision for the annual appointment of a committee of shareholders to inspect the accounts.[68] Forrester[69] describes the accounting and control system of a canal company from 1768 to 1816 and that description includes sureties or guarantees provided by officers combined with an audit at the end of an officer's term of employment. Ma and Morris examined the records of Australian and British joint stock banks prior to the 1844 Companies Act and found that those banks "had their accounts audited by the directors or provided for auditors to be appointed from the general body of shareholders."[70]

U.S. corporations also used committees of auditors.[71] The Bank of Commerce in New York in the annual report dated May 1850 reports:

[64] *Id.* at 34.

[65] 2 Scott, *supra* note 8, at 10.

[66] See Willan, *supra* note 49, at 23. The history of the Russia Company before 1666 has to be traced from sources outside the company, since like many of the other companies the Russia Company's own records were destroyed in the great fire of London. Consequently, the statement above is based on outside sources and it is possible that the accounts were audited before 1580.

[67] Henry Stevens, The Dawn of British Trade to the East Indies 107 (1967). In 1621 the East India Company appointed two paid officials (later one) with the title "Auditor." Sir William Foster, The East India House 9 (1924). This practice appears to be unusual.

[68] Armand B. DuBois, The English Business Company after the Bubble Act 1720–1800, at 300 (1938).

[69] D. A. R. Forrester, Early Canal Company Accounts: Financial and Accounting Aspects of the Forth and Clyde Navigation, 1768–1816, 10 Accounting & Bus. Res. 109 (1980).

[70] Ronald Ma & Richard D. Morris, Disclosure Practices of British and Australian Banks in the Nineteenth Century 26 (July 1980) (unpublished paper, Univ. New South Wales).

[71] See Gary J. Previts & Barbara D. Merino, A History of Accounting in America 3–4 (1979), who report that the Massachusetts Bay Company in 1629 used audit committees and that this practice continued through the 1870s.

"The usual quarterly examinations, with the usual satisfactory results have been made during the year by Committees of the Boards of Directors." The 1781 Charter of the Bank of North America specified that at every quarterly meeting of the board two directors are chosen to inspect the books every day. The Philadelphia National Bank, chartered in 1804, appointed a committee of directors to examine the books (quarterly) and "audit the bank's assets monthly."[72] Also, the Union Canal Company of Pennsylvania issued annual reports from 1824 through 1847 indicating that the accounts were examined by a committee of directors.

   3. *The Auditors' Incentives to Be Independent.*   Mechanisms are observed that provide auditors with incentives to maintain their quality. In the early joint stock companies and regulated companies a merchant's reputation affected the probability that he would be elected a director or auditor or even admitted to the company. Voting on such matters tended to be on a one-man-one-vote basis.[73] Yet some merchants appear as members of the committees or courts of directors (assistants) of several companies.[74] Further, merchants frequently appear on the audit committees of several companies in that period.[75]

   The survival of the committee of auditors for six hundred years strongly suggests it was an efficient monitoring device. Yet legal historians argue that the method was inefficient on the basis of the frauds which occurred in the seventeenth to nineteenth centuries and that government regulation was necessary for effective control.[76] However, they only consider the benefits of removing observed abuses. Ignored are the costs of regulation to remove these abuses. They fail to consider the number of firms in which the contracting and monitoring system between managers and shareholders worked. Scott makes this point in replying to Adam Smith's criticisms of joint-stock companies:

In fact, while the methods of control and of internal organization were far from perfect, they were much better than might have been expected, considering the

---

   [72] Nicholas B. Wainwright, History of the Philadelphia National Bank 20 (1953).

   [73] See Croft, *supra* note 27, at 78.

   [74] Of the eighteen directors named in the 1605 Charter of the Levant Company, at least nine also served on the committees of the East India Company and at least three also served as directors of the Spanish Company. Epstein, *supra* note 38, at 165; Croft, *supra* note 27, at 3, 4, and 247; and Stevens *supra* note 67, at 6, 7, 12, 13, 121, 179, and 237.

   [75] For example, William Harryson served on audit committees of the East India Company in 1600 and 1601 while he was a director of that company, Stevens, *supra* note 67, at 107, 156, and 166, and on an audit committee of the Spanish Company in 1605 while an ordinary member. Croft, *supra* note 27, at 37. Nicholas Lyng and Richard Wyche also served on audit committees of both the East India Company and the Spanish Company in 1601 and 1605, respectively. *Id.,* and Stevens, *supra* note 67, at 156.

   [76] Edwin M. Dodd, American Business Corporations until 1860 (1954), at 291–307.

times and how undeveloped the joint stock system was in the seventeenth century. Despite some instances of fraud, carelessness and profligacy on the part of agents abroad, numerous instances can be quoted of a remarkable devotion to duty, while amongst the directors or assistants there was a large-hearted disinterestedness, united to a careful supervision of business, which is highly commendable. *It is noteworthy that out of the great number of companies, whose officers have been investigated in this work, the allegations of fraudulent management are comparatively rare.*[77]

## D.    Observations on the Development of Auditing

The evidence suggests the practice of having a committee of auditors was not imposed on the merchant guilds, regulated companies, or joint-stock companies by law. The auditing was voluntary. Typically there were no references to auditing in the charters of the guilds, regulated companies, and joint-stock companies examined. The auditing was by order of the general court (meeting of members of shareholders) or the court of assistants (directors' meeting). Thus, when the U.K. Companies Act of 1844 required directors to keep accounts and required those accounts to be audited by persons other than the directors (or their clerks),[78] Parliament was merely incorporating into the law a version of a practice that had existed for six hundred years.

The hypothesis advanced here is that committees of auditors survived because they were an efficient method of monitoring contracts between managers and those supplying capital. However, auditing practice was not constant over time, it changed as the business corporations changed. Two dimensions in which auditing changed are the composition and relative size of audit committees. For example, the committee of auditors of the regulated and first joint-stock companies (for example, the Spanish Company and the East India Company) included assistants (that is, directors). By 1844 this was a committee of shareholders. The inclusion of assistants on the early audit committees is, at first glance, puzzling. Presumably, the assistants were managers whose actions the auditors were to monitor.

The answer appears to be that assistants were not exactly equivalent to the directors of a modern corporation. In both the Spanish Company in 1605 and the East India Company in 1600 the ratio of directors (assistants) to shareholders (ordinary members) was very large by today's standards. The Spanish Company had sixty-one assistants out of a total membership of 310.[79] The East India Company had twenty-four members of the gen-

---

[77] See Scott, *supra* note 8, at 451–52. Emphasis added.

[78] Ananias C. Littleton, Accounting Evolution to 1900, at 289 (1933).

[79] See Croft, *supra* note 27, at xxxvi. The extra assistant beyond the multiple of twelve was the secretary.

eral committee (assistants) out of a total membership of 103.[80] Given these ratios of assistants to members (one to five and one to four), the assistants were likely to be representative of the members' interests, and the large number of assistants would make it more difficult for the assistants to collude against the membership. Hence, explicit monitoring of the assistants in general would not be efficient. On the other hand, officials with direct personal control over resources (for example, the treasurer) would, in the absence of monitoring, be able to convert resources to their own personal use. Consequently, it is not surprising that the treasurer was audited by assistants and that assistants were not audited by members, in general.

If the inclusion of assistants on the audit committees of the early joint-stock companies can be explained by the high ratio of assistants or directors to shareholders, the movement to a shareholder committee can be explained by the large drop in that ratio. In the eighteenth and nineteenth centuries, despite the Bubble Act of 1720, use of the joint stock form of organization grew, and the number of directors for each company and the ratio of directors (assistants) to shareholders (members) dropped substantially, particularly in the early eighteenth century.[81] This suggests that it would be less costly for the directors to collude against the shareholders. Consequently, there was the increased tendency to use committees of shareholders, not directors, to audit the accounts.

While the evidence suggests that from the time of the guilds to the mid-nineteenth century auditing was used as a monitoring device for contracts between managers and equity holders, we found no direct evidence that audits were used to monitor debt contracts in that period. The secondary sources examined contained no references to the auditors' reporting to debt holders or that the debt holders even received the audit reports. We would not expect such evidence for the guilds, because before the Reformation loans were discouraged by the church.[82] However, in the eigh-

---

[80] Stevens, *supra* note 67, at 58–63.

[81] Two observations support this contention. First, the increase in the capitalization of joint stock companies suggests a greater number of shareholders. In 1560 the capital of the only two joint stock companies in existence was less than £10,000 or .013 percent of the national wealth of the United Kingdom. Scott, *supra* note 8, at 439. By 1695 the capital of joint stock companies was 1.3 percent of national wealth and by 1720 it was 13 percent. Scott, *supra* note 8, at 439. Also, the membership of the regulated companies and the large majority of the members of the first joint-stock companies were merchants. However, in the seventeenth and eighteenth centuries nonmerchants became the suppliers of capital to the joint-stock companies, Scott, *supra* note 8, at 440–43, suggesting an increase in the number of shareholders. Second, the average number of directors of a company appears to have *dropped* during this period. Dodd, *supra* note 76, at 291, describes the directors in the eighteenth century as a "small group."

[82] Scott, *supra* note 8, at 1.

teenth century companies issued numerous bonds.[83] Company charters under which the bonds were issued used accounting numbers (for example, profits) to restrict dividends[84] and one would expect such constraints to be monitored. There is some indirect evidence to suggest that auditing was used to monitor debt contracts. It is possible that accounts were available to bondholders and the audit by the committee of shareholders served to control the shareholder-bondholder conflict of interest.

The proposition that the accounts were available to debt holders is strengthened by the inclusion in the 1862 U.K. Companies Act of a requirement that a statement of mortgages be available to creditors.[85]

## III.   The Development of the Professional Audit Firm

In this section we provide brief descriptions of the development of the independent professional audit firm for the United Kingdom and the United States. Then explanations for the firms' development are investigated.

### A.   Development in the United Kingdom

The U.K. company acts from 1844 to 1900 did not require outside auditors. The 1844–45 acts required the directors to keep accounts and required those accounts to be audited by persons other than the directors or their clerks.[86] Further, the auditors were required to be shareholders. The 1856 act dropped the compulsory audit requirement.[87] The 1862 act included an optional model set of articles that provided for audits. While the provision did not require the auditor to be a shareholder, it also did not require the auditor be a professional firm. A series of miscellaneous acts (Railway Companies Act, 1867–68; Banking Companies, 1879; Water Companies, 1871) required audits, but again not by outsiders.[88] The 1900 Companies Act reestablished compulsory audits. However, by this time "the accounts of most of them (public companies) were not only audited but were in fact audited by chartered accountants. Indeed, practice has generally outrun legal minima."[89] Chartered accountants are professional accountants accredited by professional accounting societies.

---

[83] Dodd, *supra* note 76, at 112.

[84] Ross L. Watts & Jerold L. Zimmerman, The Demand for and Supply of Accounting Theories: The Market for Excuses, 54 Accounting Rev. 273, 277–78 (1979).

[85] Lawrence R. Dicksee, Auditing 159–66 (1892).

[86] Littleton, *supra* note 78, at 289.

[87] Bishop C. Hunt, Auditor Independence, 59 J. Accountancy 453 (1935).

[88] Francis W. Pixley, Auditors (1881); and also Dicksee, *supra* note 85.

[89] Hunt, *supra* note 87, at 454.

The substitution of professional auditors for the amateur shareholder audit committees appears to have occurred very rapidly. There is no evidence of professional firms' being appointed auditors in 1844. However, it is likely that professional firms were employed to assist the audit committee because the 1845 Companies Act allowed auditors to employ outside experts at company expense.[90] The omission of the requirement that the auditor be a shareholder from the optional articles in the 1862 act suggests pressure to appoint professionals directly. By 1881 there was a tendency to employ professional accountants directly as auditors and that tendency appears to be connected to the floating of new securities (as would be expected if an independent audit reduces the agency costs of promoters). Pixley reports the tendency thus: "[N]early all the prospectuses of new Companies now include among their officers the names of professional Accountants as their Auditors, while the older Companies are gradually replacing the Shareholders' Auditor by a professional one."[91]

## B.   Development in the United States

The 1933 Securities Act required that corporations subject to the act have audits by independent or certified public accountants. However, by the 1920s most companies listed on the New York Stock Exchange (NYSE) were already audited by professional auditors. Benston[92] reports that 82 percent of NYSE companies had professional auditors by 1926.

The substitution of professional auditors for shareholder auditors occurred later in the United States than in the United Kingdom. As noted, by 1900 most traded U.K. companies were audited by chartered accountants, but in that year only a minority of U.S.-listed companies were audited by professional auditors. A random sample of fifty-one companies listed on the NYSE in 1900 whose annual reports are in the Harvard Baker Library or the Columbia University Library revealed only eleven companies with professional auditors.[93]

Many of the U.S. audit firms existing in 1900 were started by British

---

[90] Littleton, *supra* note 78, at 296.

[91] Pixley, *supra* note 88, at 165.

[92] George J. Benston, The Effectiveness and Effects of the SEC's Accounting Disclosure Requirements, in Economic Policy and the Regulation of Corporate Securities 519 (Henry G. Manne ed. 1969).

[93] Professor David Fehr of the Harvard Business School provided us with the findings from the annual reports in the Baker Library.

chartered accountants who came to the United States to audit American companies selling securities in London.[94]

## C.  Explanations for the Professional Firm

The substitution of professional auditors for shareholder auditors occurred in both Britain and America in periods when a professional auditor was not required by law. This suggests that the substitution was the result of market forces. There were two major market developments in the period 1844–1900 in the United Kingdom that can explain the shift from shareholder to professional auditors: (1) an increase in the demand for audits and (2) the introduction of a low-cost mechanism for certifying auditor competence and independence.

1. *An Increase in the Demand for Audits.*   The demand for audits in the United Kingdom in the 1860s and 1870s increased because the complexity of the accounts, the legal liability of directors, and the size and number of corporations all increased. Accounts became more complex in the latter half of the 1800s in the United Kingdom because of government regulation of railroads and utilities and taxation that created the incentive to write off a portion of the capital stock each year.[95] In the same period courts began holding directors liable if the company was not using more complex accounting procedures (for example, provisions for doubtful accounts).[96]

Another important factor shifting the demand for auditing was the tremendous expansion in the number of companies. The nominal value of listed securities on the London Stock Exchange increased fourteenfold between 1853 and 1893.[97] While most of that increase was caused by the enormous growth in foreign securities traded (thirty-eight-fold increase), private domestic security values increased fourfold.

The increased complexity encouraged specialization in auditing and hence the growth of professional firms. The growth in the scale of the capital markets increased the fixed cost of an auditor's establishing a reputation that would serve as a bond for the auditor's independence. This led to the development of large professional audit firms.

2. *The Introduction of a Low-Cost Mechanism for Accrediting Auditors.*   The first professional society of accountants was formed in Scot-

---

[94] C. A. Moyer, Early Developments in American Auditing, 26 Acc. Rev. 3 (1951); John L. Carey, The Rise of the Accounting Profession 27 (Amer. Inst. of Cert. Pub. Accnts. 1969–70); Chester W. DeMond, Price, Waterhouse & Co. in America 10–12 (1951).

[95] Watts and Zimmerman, *supra* note 84, at 290–95.

[96] Dicksee, *supra* note 85, at 251–53; and Littleton, *supra* note 78, at 309.

[97] Edward V. Morgan & William A. Thomas, The Stock Exchange (1962), at Table 5.

land in 1854 and in England in 1870.[98] These professional societies arose to provide information on the accountant's reputation, not in auditing, but rather in bankruptcies.[99] The first bankruptcy statutes that allowed professional accountants to be appointed trustee were enacted in Scotland in 1772 and in England in 1825. Later English acts expanded the accountants' role and created a demand for information to discriminate among accountants. The 1869 act so increased the number of unskilled people competing for appointment as trustee that one judge remarked in 1875: "The whole affairs in bankruptcy have been handed over to an ignorant set of men called accountants, which was one of the greatest abuses ever introduced into the law."[100] The professional societies established brand names, thereby providing information on their members' competence and integrity. Brand names were established by setting and monitoring standards of conduct, examinations for admission,[101] and by adopting the title "chartered accountant." Once the mechanism of a professional society certifying quality in bankruptcies was established, it could be used to certify quality and independence in audits at low incremental cost.[102]

The rapid substitution of professional, "independent" auditors for lay shareholder auditors in the United Kingdom during the 1844–1900 period can be explained by (1) an outward shift in the demand curve for audits and professional auditing being a declining average cost industry[103] or (2) a downward shift in the supply curve because the start-up costs for accrediting accountants were already incurred in certifying accountants for bankruptcy work.

The development of professional audit firms being later in the United States than in the United Kingdom is consistent with the preceding explanation for the substitution of professional auditors. The United States and United Kingdom had similar rates of growth in their capital markets over

---

[98] Richard Brown, A History of Accounting and Accountants (1905), pt. 2, chs. 2–4, at 208, 235, 237.

[99] Littleton, *supra* note 78, at 271–84. The original Scottish petition for incorporation did not even list auditing as one of the public accountants' functions (Brown, *supra* note 98, at 207–08.

[100] Littleton, *supra* note 78, at 283.

[101] Nicholas A. H. Stacey, English Accountancy 21 (1954).

[102] A non–mutually exclusive hypothesis for the establishment of professional societies is that they were designed to restrain trade. However, it does not appear that this motivation was important, because the professional societies in the United Kingdom did not take any action to restrict entry until 1893, when the Institute of Chartered Accountants in England and Wales had a bill introduced in Parliament "to restrain all persons from practising who are not Chartered Accountants." Brown, *supra* note 98, at 243. Moreover, this bill was withdrawn and other societies competed with the Institute for the next seventy years.

[103] Stigler, *supra* note 51.

the years from 1853 to 1903,[104] but the U.K. capital market was much larger than in the United States in 1853. Thus, the absolute increase in the scale of the capital markets and in the demand for auditing in the period 1844–1900 was larger in the United Kingdom. In addition, the United States did not experience the same reduction in start-up costs for accrediting professional accountants.

Unlike the United Kingdom, the United States did not experience the change in bankruptcy laws that gave creditors the powers to choose accountants as trustees. Hence, there was less demand in the United States for information on an accountant's reputation for handling bankruptcies, and consequently, professional accounting societies were not established to accredit accountants.

The first U.S. professional accountants society was formed in 1887,[105] thirty-three years after the first Scottish society and seventeen years after the first English society. At that time most of the accountants' work in the United States involved audits and investigations, and little or none involved liquidations and bankruptcies.[106] The start-up costs of the American societies were borne largely by British accountants who came to America to audit firms raising capital in London and stayed to start their own firms. Evidence of the value of the British auditors' brand name capital is provided by the British auditors' urging the American society to adopt the title "certified public accountant" to designate their members instead of the title of "chartered accountant."[107]

The first U.S. society began accrediting members in 1896, when the first certified public accountants law in New York state was passed. Following this law, nonaccredited auditors were rapidly replaced, so that, as noted, by the 1920s most NYSE companies were audited by professionals.

The change from shareholder auditors to professional auditors was not a switch from nonindependent to independent auditors, if independence is interpreted as the likelihood that an auditor will report a discovered breach of contract. An important incentive to be independent, the effect of failure to report a breach on the auditor's reputation and hence his future business, existed for both the shareholder and professional auditors.

---

[104] Ross L. Watts & Jerold L. Zimmerman, The Markets for Independence and Independent Auditors 30 (June 1981) (unpublished manuscript, Univ. Rochester, Grad. School Management).

[105] Carey, *supra* note 94, at 36–39.

[106] Brown, *supra* note 98, at 278–79.

[107] James D. Edwards, The Antecedents of American Public Accounting (1956), reprinted in Contemporary Studies in the Evolution of Accounting Thought 53 (Michael Chatfield ed. 1968).

An interesting example of the importance of reputation, or brand name, to the auditor is provided by the expansion of Price, Waterhouse and Company to the United States in the late nineteenth century. Price, Waterhouse and Company did not allow the (nonpartner) representatives they sent to America (Jones and Caesar) to use the firm name for fear of damage to Price Waterhouse's reputation.[108]

## IV. CONCLUSIONS

The survival of the bonding and auditing practices from the Ipswich merchant guild in 1200 to the British joint stock banks of 1836 and U.S. banks and canal companies of the mid-nineteenth century is consistent with the existence of agency problems and the use of bonding and monitoring devices to reduce agency costs.[109] Moreover, the size and composition of audit committees are observed evolving in response to changes in the size of the firm and the nature of the agency costs faced by the contracting parties.

The pervasiveness of voluntary (costly) auditing in the precursors to the modern corporation is consistent also with auditors' developing quality-assuring devices—in particular, mechanisms that increase the probability the auditor will report a breach in a contract he is to monitor (that is, be independent). We observe mechanisms being devised that supplied the incentives for auditors to maintain their independence in the guild and regulated companies (committees and penalties, including loss of reputation), in the joint stock companies, and finally in the development of professional societies.

Overall, the evidence suggests that the existence of the independent auditor is not the direct result of government fiat. The appearance of the professional independent auditor was encouraged by changes in U.K. bankruptcy laws, but the United States' evidence suggests that even without those bankruptcy laws, economies of scale in auditing would have led to the development of the professional independent auditor.

---

[108] DeMond, *supra* note 94, at 16.

[109] T. A. Lee, The Historical Development of Internal Control from the Earliest Times to the End of the Seventeenth Century, 9 J. Accounting Res. 150 (1971), traces the use of audits in government back to ancient times.

# AN ANALYSIS OF THE PRINCIPAL-AGENT PROBLEM

## By Sanford J. Grossman and Oliver D. Hart[1]

Most analyses of the principal-agent problem assume that the principal chooses an incentive scheme to maximize expected utility subject to the agent's utility being at a stationary point. An important paper of Mirrlees has shown that this approach is generally invalid. We present an alternative procedure. If the agent's preferences over income lotteries are independent of action, we show that the optimal way of implementing an action by the agent can be found by solving a convex programming problem. We use this to characterize the optimal incentive scheme and to analyze the determinants of the seriousness of an incentive problem.

## 1. INTRODUCTION

IT HAS BEEN RECOGNIZED for some time that, in the presence of moral hazard, market allocations under uncertainty will not be unconstrained Pareto optimal (see Arrow [1], Pauly [13]). It is only relatively recently, however, that economists have begun to undertake a systematic analysis of the properties of the second-best allocations which will arise under these conditions. Much of this analysis has been concerned with what has become known as the principal-agent problem. Consider two individuals who operate in an uncertain environment and for whom risk sharing is desirable. Suppose that one of the individuals (known as the agent) is to take an action which the other individual (known as the principal) cannot observe. Assume that this action affects the total amount of consumption or money which is available to be divided between the two individuals. In general, the action which is optimal for the agent will depend on the extent of risk sharing between the principal and the agent. The question is: What is the optimal degree of risk sharing, given this dependence?

Particular applications of the principal-agent problem have been made to the case of an insurer who cannot observe the level of care taken by the person being insured; to the case of a landlord who cannot observe the input decision of a tenant farmer (sharecropping); and to the case of an owner of a firm who cannot observe the effort level of a manager or worker.[2]

Although considerable progress has been made in the recent literature towards understanding and solving the principal-agent problem (see, in particular, Harris and Raviv [6], Holmstrom [7], Mirrlees [10, 11, 12], Shavell [19, 20], as well as the other references in footnote 2), the mathematical approach which has been adopted in most of this literature is unsatisfactory. The procedure usually followed is to suppose that the principal chooses the risk-sharing contract, or incentive scheme, to maximize his expected utility subject to the constraints that

[2] These and other applications are discussed in a number of recent papers. See, for example, Harris and Raviv [6], Holmstrom [7], Mirrlees [10, 11, 12], Radner [15], Ross [17], Rubinstein and Yaari [18], Shavell [19, 20], Spence and Zeckhauser [21], Stiglitz [22], and Zeckhauser [24].

7

For a given $\underset{\sim}{I}$ the agent strictly prefers lower actions

FIGURE 1.

(a) the agent's expected utility is no lower than some pre-specified level; (b) the agent's utility is at a stationary point, i.e., the agent satisfies his first-order conditions with respect to the choice of action. That is, the agent's second-order conditions (and the condition that the agent should be at a global rather than a local maximum) are ignored. Mirrlees [10], however, in an important paper, has shown that this procedure is generally invalid unless, at the optimum, the solution to the agent's maximum problem is unique. In the absence of uniqueness (and it is difficult to guarantee uniqueness in advance), the first-order conditions derived by the above procedure are not even necessary conditions for the optimality of the risk-sharing contract.[3]

---

[3] The reason for this can be seen quite easily in Figure 1 (we are grateful to Andreu Mas-Colell for suggesting the use of this figure). On the horizontal axis, $I$ represents the agent's incentive scheme and on the vertical axis $a$ represents the agent's action. The curve $ABCDE$ is the locus of pairs of actions and incentive schemes which satisfy the agent's first order conditions, i.e., given $I$ the agent's utility is at a stationary point. Of these points, only those lying on the segments $AB$ and $DE$ represent global maxima for the agent, e.g. given the incentive scheme $\underset{\sim}{I}$ the agent's optimal action is at $p_1$, not at $p_2$ or $p_3$. Indifference curves—in terms of $a$ and $I$—are drawn for the principal ($C$ is on a higher curve than $B$). The true feasible set for the principal are the segments $AB$ and $DE$ and the optimal outcome for the principal is therefore $B$. However, $B$ does not satisfy the first order conditions of the problem: maximize the principal's utility subject to $(a, I)$ lying on $ABCDE$, i.e., subject to $(a, I)$ satisfying the agent's first order conditions (the solution to this problem is at $C$). In other words, $B$ does not satisfy the necessary conditions for optimality of the problem which has been studied in much of the literature. Note finally that perturbing Figure 1 slightly does not alter this conclusion.

The purpose of this paper is to develop a method for analyzing the principal-agent problem which avoids the difficulties of the "first-order condition" approach.[4] Our approach is to break the principal's problem up into a computation of the costs and benefits of the different actions taken by the agent. For each action, we consider the incentive scheme which minimizes the (expected) cost of getting the agent to choose that action. We show that, under the assumption that the agent's preferences over income lotteries are independent of the action he takes, this cost minimization problem is a fairly straightforward convex programming problem. An analysis of these convex problems as the agent's action varies yields a number of results about the form of the optimal incentive scheme. We will also be able to analyze what factors determine how serious a particular incentive problem is; i.e., how great the loss is to the principal from having to operate in a second-best situation where the agent's action cannot be observed relative to a first-best situation where it can be observed.

The assumption that the agent's preferences over income lotteries are independent of action is a strong one. Yet it seems a natural starting point for an analysis of the principal-agent problem. Special cases of this assumption occur when the agent's utility function is additively or multiplicatively separable in action and reward. One or other of these cases is typically assumed in most of the literature. In Section 6 we discuss briefly the prospects for the non-independence case.

In addition to providing greater rigor, the costs versus benefits approach also provides a clear separation of the two distinct roles the agent's output plays in the principal-agent problem. On the one hand, the agent's output contributes positively to the principal's consumption, so the principal desires a high output. On the other hand, the agent's output is a signal to the principal about the agent's level of effort. This informational role may be in conflict with the consumption role. For example, there may be a moderate output level which is achieved when the agent takes low effort levels and never occurs at other effort levels. If the agent is penalized whenever this moderate output occurs, then he is discouraged from taking these low effort actions. However, there may be lower output levels which have some chance of occurring regardless of the agent's action. To encourage the agent to take high effort levels, it is then optimal to pay the agent more in low output states than in moderate output states, even though the principal prefers moderate output levels to low output levels.

The dual role of output makes it difficult to obtain conditions which ensure even elementary properties of the incentive scheme, such as monotonicity. In Section 3, sufficient conditions for monotonicity are given. It is also shown in this section that a monotone likelihood ratio condition, which the "first-order condition" approach suggests is a guarantee of monotonicity, must be strengthened once we take into account the possibility that the agent's action is not unique at the optimal incentive scheme.

The paper is organized as follows. In Section 2, we show how the principal's optimization problem can be decomposed into a costs versus benefits problem.

---

[4]Mirrlees [12] has identified a class of cases where the "first-order condition" approach *is* valid. We will consider this class in Section 3.

In Section 3, we use our approach to analyze the monotonicity and progressivity of the optimal incentive scheme. In Section 4, we give a simple algorithm for computing an optimal incentive scheme when there are only two outcomes associated with the agent's actions. In Section 5, we analyze the effects of risk aversion and information quality on the incentive problem. Finally, in Section 6 we consider some extensions of the analysis.

## 2. STATEMENT OF THE PROBLEM

The application of the principal-agent problem that we will consider is to the case of the owner of a firm who delegates the running of the firm to a manager. The owner is the principal and the manager the agent. The owner is assumed not to be able to monitor the manager's actions. The owner does, however, observe the outcome of these actions, which we will take to be the firm's profit. It is assumed that the firm's profit depends on the manager's actions, but also on other factors which are outside the manager's control—we model these as a random component. Thus, in particular, if the firm does well, it will not generally be clear to the owner whether this is because the manager has worked well or whether it is because he has been lucky.[5]

We will simplify matters by assuming that there are only finitely many possible gross profit levels for the firm, denoted $q_1, \ldots, q_n$, where $q_1 < q_2 < \cdots < q_n$. We will assume that the principal is interested only in the firm's net profit, i.e. gross profit minus the payment to the manager. We will also assume that the principal is risk neutral—our methods of analysis can, however, be applied to the case where the principal is risk averse (see Remark 3 and Section 6).

Let $A$ be the set of actions available to the manager. We will assume that $A$ is a non-empty, compact subset of a finite dimensional Euclidean space. Let $S = \{x \in R^n \mid x \geq 0, \sum_{i=1}^{n} x_i = 1\}$. We assume that there is a continuous function $\pi : A \to S$, where $\pi(a) = (\pi_1(a), \ldots, \pi_n(a))$ gives the probabilities of the $n$ outcomes $q_1, \ldots, q_n$ if action $a$ is selected. It is assumed that, when the agent chooses $a \in A$, he knows the probability function $\pi$ but not the outcome which will result from his action. We assume that the agent has a von Neumann–Morgenstern utility function $U(a, I)$ which depends both on his action $a$ and his remuneration $I$ from the principal. We include $a$ as an argument in order to capture the idea that the agent dislikes working hard, taking care, etc.

The crucial assumption that we will make about the form of $U(a, I)$ is:

ASSUMPTION A1: $U(a, I)$ can be written as $G(a) + K(a)V(I)$, where (1) $V$ is a real-valued, continuous, strictly increasing, concave function defined on some open interval $\mathcal{I} = (\underline{I}, \infty)$ of the real line; (2) $\lim_{I \to \underline{I}} V(I) = -\infty$; (3) $G, K$ are

---

[5]The assumption that the principal cannot monitor the agent's actions at all may in some cases be rather extreme. For a discussion of the implications of the existence of imperfect monitoring opportunities, see Harris and Raviv [6], Holmstrom [7] and Shavell [19, 20]. See also Remark 4 in Section 2.

real-valued, continuous functions defined on $A$ and $K$ is strictly positive; (4) for all $a_1, a_2 \in A$ and $I, \hat{I} \in \mathcal{I}$, $G(a_1) + K(a_1)V(I) \geq G(a_2) + K(a_2)V(I) \Rightarrow G(a_1) + K(a_1)V(\hat{I}) \geq G(a_2) + K(a_2)V(\hat{I})$.

In the above, we allow for the case $\underline{I} = -\infty$.

The main part of Assumption A1 has a simple ordinal interpretation. Assumption A1 implies that the agent's preferences over income lotteries are independent of his action (Assumption A1(1) tells us also that these preferences exhibit risk aversion). The converse can also be shown to be true: if the agent's preferences over income lotteries are independent of $a$, then $U$ can be written as $G(a) + K(a)V(I)$ for some functions $G, K, V$ (for a proof, see Keeney [8]). Note that Assumption A1 does not imply that the agent's preferences for *action* lotteries are independent of income. We will insist, however, that the agent's ranking over *perfectly certain* actions is independent of income—this is condition (4) of Assumption A1.

Note that if $K(a)$ is not constant then (2) and (4) imply that $V(I)$ must be bounded from above. Further if it is also the case that $G(a) \equiv 0$, then $V(I)$ must be non-positive everywhere.

Two special cases of Assumption A1 occur when $K(a) = $ constant, i.e. $U$ is additively separable in $a$ and $I$, and when $G(a) = 0$, i.e. $U$ is multiplicatively separable in $a$ and $I$. In these cases the agent's preferences over action lotteries *are* independent of income, as well as preferences over income lotteries being independent of action.[6]

An interesting special case of multiplicative separability is when $V(I) = -e^{-kI}$, $K(a) = e^{ka}$ and $A$ is a subset of the real line. Then $U(a, I) = -e^{-k(I-a)}$; i.e., effort appears just as negative income.

In the "first-best" situation where the principal can observe $a$, it is optimal for him to pay the agent according to the action he chooses. Let $\overline{U}$ be the agent's reservation price, i.e. the expected level of utility he can achieve by working elsewhere, and let $\mathcal{U} = V(\mathcal{I}) = \{v \,|\, v = V(I) \text{ for some } I \in \mathcal{I}\}$. We make the following assumption.

ASSUMPTION A2: $[\overline{U} - G(a)]/K(a) \in \mathcal{U}$ for all $a \in A$.

DEFINITION: Let $C_{FB} : A \to R$ be defined by $C_{FB}(a) = h([\overline{U} - G(a)]/K(a))$, where $h \equiv V^{-1}$.

Here $C_{FB}$ stands for first-best cost. $C_{FB}(a)$ is simply the agent's reservation price for picking action $a$. To get the agent to pick $a \in A$ in the first-best

---

[6] The converse is also true: if preferences over action lotteries are independent of income as well as preferences over income lotteries being independent of action, then $U$ is additively or multiplicatively separable (see Keeney [8] or Pollak [14]).

situation, the principal will offer him the following contract: I will pay you $C_{FB}(a)$ if you choose $a$ and $\tilde{I}$ otherwise, where $\tilde{I}$ is very close to $\underline{I}$.

DEFINITION: Let $B : A \to R$ be defined by $B(a) = \sum_{i=1}^{n} \pi_i(a) q_i$. $B(a)$ is the expected benefit to the principal from getting the agent to pick $a$.

DEFINITION: A first-best optimal action is one which maximizes $B(a) - C_{FB}(a)$ on $A$.

The function $C_{FB}$ induces a complete ordering on $A : a \gtrsim a'$ if and only if $C_{FB}(a) \geq C_{FB}(a')$. For obvious reasons we will refer to actions with higher $C_{FB}(a)$'s as costlier actions. It is easy to show, in view of Assumption A1(4), that $C_{FB}(a) \geq C_{FB}(a') \Leftrightarrow G(a) + K(a)v \leq G(a') + K(a')v$ for all $v \in \mathfrak{U} \Leftrightarrow G(a) + K(a)v \leq G(a') + K(a')v$ for some $v \in \mathfrak{U}$. This in turn implies that the ordering $\gtrsim$ is independent of $\bar{U}$. In the second-best situation where $a$ is not observed by the principal, it is not possible to make the agent's remuneration depend on $a$. Instead, the principal will pay the agent according to the *outcome* of his action, i.e. according to the firm's profit. An incentive scheme is therefore an $n$-dimensional vector $I = (I_1, I_2, \ldots, I_n) \in \mathfrak{I}^n$, where $I_i$ is the agent's remuneration in the event that the firm's profit is $q_i$. Given the incentive scheme $I$, the agent will choose $a \in A$ to maximize $\sum_{i=1}^{n} \pi_i(a) U(a, I_i)$.

We will assume that the principal knows the agent's utility function $U(a, I)$, the set $A$ and the function $\pi : A \to S$. In other words, the principal is fully informed about the agent and about the firm's production possibilities. The incentive problem which we will study therefore arises entirely because the principal cannot monitor the agent's actions.[7]

The principal's problem can be described as follows. Let $F$ be the set of pairs of incentive schemes $I^*$ and actions $a^*$ such that, under $I^*$, the agent will be willing to work for the principal and will find it optimal to choose $a^*$, i.e. $\max_{a \in A} \sum_{i=1}^{n} \pi_i(a) U(a, I_i^*) = \sum_{i=1}^{n} \pi_i(a^*) U(a^*, I_i^*) \geq \bar{U}$. Then the principal chooses $(I, a) \in F$ to maximize $\sum_{i=1}^{n} \pi_i(a)(q_i - I_i)$. It simplifies matter considerably if we break this problem up into two parts. We consider first, given that the principal wishes to implement $a^*$, the least cost way of achieving this. We then consider which $a^*$ should be implemented. Thus, to begin, suppose that the principal wishes the agent to pick a particular action $a^* \in A$. To find the least (expected) cost way of achieving this, the principal must solve the following

---

[7] This distinguishes our study from the literature on incentive compatibility; see, e.g., the recent *Review of Economic Studies* symposium [16]. The incentive compatibility literature has been concerned with incentive problems arising from differences in information between individuals rather than with those arising from monitoring problems. In cases of differential information, there is a role for an exchange of information through messages, whereas in the model we study messages would serve no purpose.

problem:

(2.1)     Choose $I_1, \ldots, I_n$ to minimize $\sum_{i=1}^{n} \pi_i(a^*)I_i$

        subject to $\sum_{i=1}^{n} \pi_i(a^*)U(a^*, I_i) \geq \bar{U},$

$$\sum_{i=1}^{n} \pi_i(a^*)U(a^*, I_i) \geq \sum_{i=1}^{n} \pi_i(a)U(a, I_i) \qquad \text{for all} \quad a \in A,$$

$$I_i \in \mathcal{I} \qquad \text{for all } i.$$

This problem can be simplified considerably in view of Assumption A1. It will be convenient to regard $v_1 = V(I_1), \ldots, v_n = V(I_n)$ as the principal's control variables. Recall that $\mathcal{U} = V(\mathcal{I}) = \{v \mid v = V(I) \text{ for some } I \in \mathcal{I}\}$. By Assumption A1, $\mathcal{U}$ is an interval of the real line $(-\infty, \bar{v})$. Thus we may rewrite (2.1) as follows:

(2.2)     Choose $v_1, \ldots, v_n$ to minimize $\sum_{i=1}^{n} \pi_i(a^*)h(v_i)$

        subject to $G(a^*) + K(a^*)\left( \sum_{i=1}^{n} \pi_i(a^*)v_i \right) \geq G(a) + K(a)\left( \sum_{i=1}^{n} \pi_i(a)v_i \right)$

$$\text{for all} \quad a \in A,$$

$$G(a^*) + K(a^*)\left( \sum_{k=1}^{n} \pi_i(a^*)v_i \right) \geq \bar{U},$$

$$v_i \in \mathcal{U} \qquad \text{for all } i,$$

where $h \equiv V^{-1}$.

The important point to realize is that the constraints in (2.2) are linear in the $v_i$'s. Furthermore, $V$ concave implies $h$ convex, and so the objective function is convex in the $v_i$'s. Thus (2.2) is a rather simple optimization problem: minimize a convex function subject to (a possibly infinite number of) linear constraints. In particular, when $A$ is a finite set, the Kuhn–Tucker theorem yields necessary and sufficient conditions for optimality. These will be analyzed later.

It is important to realize that, in the absence of Assumption A1, it is not generally possible to convert (2.1) into a convex problem in this way.

DEFINITION: If $I = (I_1, \ldots, I_n)$ satisfies the constraints in (2.1) or $v = (v_1, \ldots, v_n)$ satisfies the constraints in (2.2), we will say that $I$ or $v$ *implements* action $a^*$. (We are assuming here that if the agent is indifferent between two actions, he will choose the one preferred by the principal.)

Consider the set of $v$'s which implement $a^*$. For some $a^*$, this set may be empty, in which case action $a^*$ cannot be implemented by the principal at any cost. If the set is non-empty, then, since $h$ is convex,

$$\sum_{i=1}^{n} \pi_i(a^*)h(v_i) \geq h\left(\sum_{i=1}^{n} \pi_i(a^*)v_i\right) \geq h\left(\frac{\overline{U} - G(a^*)}{K(a^*)}\right)$$

by (2.2), and so the principal's objective function is bounded below on this set. Let $C(a^*)$ be the greatest lower bound of $\sum_{i=1}^{n}\pi_i(a^*)h(v_i)$ on this set.

DEFINITION: Let $C(a^*) = \inf\{\sum_{i=1}^{n}\pi_i(a^*)h(v_i) \mid v = (v_1, \ldots, v_n)$ implements $a^*\}$ if the constraint set in (2.2) is non-empty. In the case where the constraint set of (2.2) is empty, write $C(a^*) = \infty$. This defines the second-best cost function $C: A \rightarrow Ru\{\infty\}$.

The above constitutes the first step(s) of the principal's optimization problem: for each $a \in A$, compute $C(a)$. The second step is to choose which action to implement, i.e. to choose $a \in A$ to maximize $B(a) - C(a)$. This second problem will not generally be a convex problem. This is because even if $B(a)$ is concave in $a$, $C(a)$ will not generally be convex. Fortunately, a significant amount of information about the form of the optimal incentive scheme can be obtained by studying the first step alone.

DEFINITION: A second-best optimal action $\hat{a}$ is one which maximizes $B(a) - C(a)$ on $A$. A second-best optimal incentive scheme $\hat{I}$ is one that implements a second-best optimal action $\hat{a}$ at least expected cost, i.e. $\sum_{i=1}^{n}\pi_i(\hat{a})\hat{I}_i = C(\hat{a})$.

Note that for a second-best optimal incentive scheme to exist, the greatest lower bound in the definition of $C(a)$ must actually be achieved. In order to establish the existence of a second-best optimal action and a second-best optimal incentive scheme, we need a further assumption.

ASSUMPTION A3: For all $a \in A$ and $i = 1, \ldots, n$, $\pi_i(a) > 0$.

Since there are only finitely many possible profit levels, Assumption A3 implies that $\pi_i(a)$ is bounded away from zero. Hence Assumption A3 rules out cases studied by Mirrlees [12] in which an optimum can be approached but not achieved by imposing higher and higher penalties on the agent which occur with smaller and smaller probability if the agent chooses the right action.

PROPOSITION 1: Assume A1–A3. Then there exists a second-best optimal action $\hat{a}$ and a second-best optimal incentive scheme $\hat{I}$.

PROOF: It is helpful to split the proof up into two parts. Consider first the case where $V$ is linear. Then it is easy to see that the principal can do as well in the second-best as in the first-best where the agent can be monitored. For let $a^*$ maximize $B(a) - C_{FB}(a)$ on $A$. Let the principal offer the agent the incentive scheme $I_i = q_i - t$, where $t = B(a^*) - C_{FB}(a^*)$. Then the principal's profit will be $B(a^*) - C_{FB}(a^*)$ whatever the agent does. On the other hand, by picking $a = a^*$, the agent can obtain expected utility $\overline{U}$. Hence Proposition 1 certainly holds when $V$ is linear.

On the other hand, suppose $V$ is not linear. We show first, that, if the constraint set is nonempty for an action $a^* \in A$, then problem (2.2) has a solution, i.e. $\sum_{i=1}^n \pi_i(a^*)h(v_i)$ achieves its greatest lower bound $C(a^*)$. Note that $\sum_{i=1}^n \pi_i(a^*)v_i$ is bounded below on the constraint set of (2.2). It therefore follows from a result of Bertsekas [2] that unbounded sequences in the constraint set make $\sum_{i=1}^n \pi_i(a^*)h(v_i)$ tend to infinity (roughly because the variance of the $v_i \to \infty$ while their mean is bounded below, and $h$ is convex and nonlinear— Assumption A3 is important here). Hence, we can artificially bound the constraint set. Since the constraint set is closed, the existence of a minimum therefore follows from Weierstrass' theorem.

We show next that $C(a)$ is a lower semicontinuous function of $a$. If $A$ is finite, then any function defined on $A$ is continuous and hence lower semicontinuous. Assume therefore that $A$ is not finite. Let $(a_r)$ be a sequence of points in $A$ converging to $a$. Assume without loss of generality (w.l.o.g.) that $C(a_r) \to k$. Then, if $k = \infty$, we certainly have $C(a) \leq \lim_{r \to \infty} C(a_r)$. Suppose therefore that $k < \infty$. Let $(I_1^r, \ldots, I_n^r)$ be the solution to (2.1) when $a^* = a_r$. Then Bertsekas' result together with Assumption A3 shows that the sequence $((I_1^r, \ldots, I_n^r))$ is bounded (otherwise $C(a_r) \to \infty$). Let $(I_1, \ldots, I_n)$ be a limit point. Then clearly $(I_1, \ldots, I_n)$ implements $a$ and so $C(a) \leq \sum_{i=1}^n \pi_i(a)I_i = \lim_{r \to \infty} C(a_r)$. This proves lower semicontinuity.

Given that $C(a)$ is lower semicontinuous and $A$ is compact, it follows from Weierstrass' theorem that $\max_{a \in A}[B(a) - C(a)]$ has a solution, as long as $C(a)$ is finite for some $a \in A$. To prove this last part, we show that $C(a^*) = C_{FB}(a^*)$ if $a^*$ minimizes $C_{FB}(a)$ on $A$. To see this, note that the $a^*$ which minimizes $C_{FB}(a)$ can be implemented by setting $I_i = C_{FB}(a^*)$ for all $i$.

We have thus established the existence of a second-best optimal action, $\hat{a}$, when $V$ is nonlinear. Since we have also shown that (2.2) has a solution as long as the constraint set is non-empty and $V$ is nonlinear, this establishes the existence of a second-best optimal incentive scheme. $Q.E.D.$

It is interesting to ask whether the constraint that the agent's expected utility be greater than or equal to $\overline{U}$ is binding at a second-best optimum. The answer is no in general, i.e. for incentive reasons it may pay the principal to choose an incentive scheme which gives the agent an expected utility in excess of $\overline{U}$. One case where this will not happen is when the agent's utility function is additively or multiplicatively separable in action and reward:

PROPOSITION 2: *Assume* A1, A2, *and either* $K(a)$ *is a constant function on* $A$ *or* $G(a) = 0$ *for all* $a \in A$. *Let* $\hat{a}$ *be a second-best optimal action and* $\hat{I}$ *a second-best optimal incentive scheme which implements* $\hat{a}$. *Then* $\sum_{i=1}^{n} \pi_i(\hat{a}) U(\hat{a}, \hat{I}_i) = \overline{U}$.

PROOF: Suppose not. Write $\hat{v}_i = V(\hat{I}_i)$. Then $G(\hat{a}) + K(\hat{a})(\sum_{i=1} \pi_i(\hat{a})\hat{v}_i) > \overline{U}$ in (2.2). But it is clear that the principal's costs can be reduced and all the constraints of (2.2) will still be satisfied if we replace $\hat{v}_i$ by $(\hat{v}_i - \epsilon)$ for all $i$ in the additively separable case and by $v_i(1 + \epsilon)$ for all $i$ in the multiplicatively separable case where $\epsilon > 0$ is small (note that in the multiplicatively separable case, it follows from (2)–(4) of Assumption A1 that $V(I) < 0$ for all $I \in \mathcal{I}$, and so $\hat{v}_i < 0$). In other words, $\hat{a}$ can be implemented at lower expected cost, which contradicts the fact that we are at a second-best optimum.               Q.E.D.

REMARK 1: The proof of Proposition 1 establishes that $C(a^*) = C_{FB}(a^*)$ if $a^*$ minimizes $C_{FB}(a)$ on $A$. This is a reflection of the fact that there is no trade-off between risk sharing and incentives when the action to be implemented is a cost-minimizing one (i.e. involves the agent in minimum "effort").

REMARK 2: In general, there may be more than one second-best optimal action and more than one second-best optimal incentive scheme. It is clear from (2.2), however, that, if $V$ is strictly concave, there is a unique second-best optimal incentive scheme which implements any particular second-best optimal action.

DEFINITION: Let $L = \max_{a \in A}(B(a) - C_{FB}(a)) - \sup_{a \in A}(B(a) - C(a))$ be the difference between the principal's expected profit in the first-best and second-best situations.

$L$ represents the loss which the principal incurs as a result of being unable to observe the agent's action (we write $\sup(B(a) - C(a))$ rather than $\max(B(a) - C(a))$ to cover cases where the assumptions of Proposition 1 do not hold). Proposition 3 shows that, while there are some special cases in which $L = 0$, in general $L > 0$.

PROPOSITION 3: *Assume* A1 *and* A2. *Then:* (1) $C(a) \geq C_{FB}(a)$ *for all* $a \in A$, *which implies that* $L \geq 0$. (2) *If* $V$ *is linear*, $L = 0$. (3) *If there exists a first-best optimal action* $a^* \in A$ *satisfying: for each* $i$, $\pi_i(a^*) > 0 \Rightarrow \pi_i(a) = 0$ *for all* $a \in A$, $a \neq a^*$, *then* $L = 0$. (4) *If* $A$ *is a finite set and there is a first-best optimal action* $a^*$ *which satisfies: for some* $i$, $\pi_i(a^*) = 0$ *and* $\pi_i(a) > 0$ *for all* $a \in A$, $a \neq a^*$, *then* $L = 0$. (5) *If there is a first-best optimal action* $a^* \in A$ *which minimizes* $C_{FB}(a)$ *on* $A$, $L = 0$. (6) *If Assumption* A3 *holds, every maximizer* $\tilde{a}$ *of* $B(a) - C_{FB}(a)$ *on* $A$ *satisfies* $C_{FB}(\tilde{a}) > \min_{a \in A} C_{FB}(a)$, *and* $V$ *is strictly concave, then* $L > 0$.

PROOF: (1) is obvious since anything which is second-best feasible is also first-best feasible. (2) follows from the first part of the proof of Proposition 1. (5) follows from the proof of Proposition 1 (see also Remark 1). (3) and (4) follow from the fact that $a^*$ can be implemented by offering the agent $I_i = C_{FB}(a^*)$ for those $i$ such that $\pi_i(a^*) > 0$ and $I$ close to $\underline{I}$ otherwise.

To prove (6), note that, if $V$ is strictly concave,

$$G(a^*) + K(a^*) \sum_{i=1}^{n} \pi_i(a^*) V(I_i) \geq \overline{U}$$

implies

$$C(a^*) = \sum_{i=1}^{n} \pi_i(a^*) h(V(I_i))$$

$$> h\left( \sum_{i=1}^{n} \pi_i(a^*) V(I_i) \right) \geq h\left( (\overline{U} - G(a^*))/K(a^*) \right)$$

$$= C_{FB}(a^*)$$

unless $I_i$ = constant with probability 1. But, since $\pi_i(a^*) > 0$ for all $i$, $I_i$ = constant with probability $1 \Rightarrow I_i$ is independent of $i$. However, in this case, the constraints of problem (2.2) imply that $C_{FB}(a)$ is minimized at $a^*$. $\quad Q.E.D.$

Most of Proposition 3 is well known. Proposition 3(2) and (6) can be understood as follows. In the first-best situation, if the agent is strictly risk averse, the principal bears all the risk and the agent bears none. In the second best situation, this is generally undesirable. For if the agent is completely protected from risk, then he has no incentive to work hard; i.e., he will choose $a \in A$ to minimize $C_{FB}(a)$. Hence the second-best situation is strictly worse from a welfare point of view than the first-best situation. The exception is when the agent is risk neutral, in which case it is optimal both from a risk sharing and an incentive point of view for him to bear all the risk, or when the first-best optimal action is cost minimizing.

In the case of Proposition 3(3) and 3(4), a scheme in which the agent is penalized very heavily if certain outcomes occur can be used to achieve the first best. This relates to results obtained in Mirrlees [12].

REMARK 3: We have assumed that the principal is risk neutral. Our analysis generalizes to the case where the principal is risk averse, however. In this case, instead of choosing $v$ to minimize $\sum \pi_i(a^*) h(v_i)$ in problem (2.2), we choose $v$ to maximize $\sum \pi_i(a^*) U_p(q_i - h(v_i))$, where $U_p$ is the principal's utility function. Note that (2.2) is still a convex problem. Although we can no longer analyze costs and benefits separately, we can, for each $a^* \in A$, define a net benefit function $\max_v \sum \pi_i(a^*) U_p(q_i - h(v_i))$. An optimal action for the principal is now one that maximizes net benefits. See also Section 6 on this.

REMARK 4: We have taken the outcomes observed by the principal to be profit levels. Our analysis generalizes, however, to the case where the outcomes are more complicated objects, such as vectors of profits, sales, etc., or to the case where profits are not observed at all but something else is (see, e.g., Mirrlees [11]). The important point to realize is that profit does not appear in the cost

minimization problem (2.1) or (2.2). Thus, if the principal observes the realizations of a signal $\tilde{\theta}$, then $I_i$ refers to the payment to the agent when $\tilde{\theta} = \theta_i$. Let $\hat{C}(a, \tilde{\theta})$ be the cost of implementing $a$ when the information structure is $\tilde{\theta}$ (e.g. if $\tilde{\theta}$ reveals $a$ exactly, then $\hat{C}(a, \tilde{\theta}) = C_{FB}(a)$). Note that if the distribution of output is generated by a production function $f(a, \tilde{w})$, such that the marginal distribution of $\tilde{w}$ is independent of the information structure, then $B(a) = Ef(a, \tilde{w}) = E[E[f(a, \tilde{w}) | \theta]]$ is independent of the information structure, given $a$. It follows that the effect of changes in the information structure is summarized by the way that $C(a, \tilde{\theta})$ changes when the information structure changes. As will be seen in Section 5, this is quite easy to analyze.

### 3. SOME CHARACTERISTICS OF OPTIMAL INCENTIVE SCHEMES

It is of interest to know whether the optimal incentive scheme is monotone increasing (i.e., whether the agent is paid more when a higher output is observed) and whether the scheme is progressive (i.e., whether the marginal benefit to the agent of increased output is decreasing in output). These questions are quite difficult to answer because of the informational role of output. As we noted in the introduction, the agent may be given a low income at intermediate levels of output in order to discourage particular effort levels. Nevertheless, some general results about the shape of optimal schemes can be established. We begin with the following lemma.

LEMMA 1: *Assume A1–A3. Let $(I_i)_{i=1}^n, (I_i')_{i=1}^n$ be incentive schemes which cause $a$ and $a'$ to be optimal choices for the agent, respectively, and minimize the respective costs (i.e. (2.1) or (2.2) is solved). Let $v_i = V(I_i)$ and $v_i' = V(I_i')$. Then, if $G(a) + K(a)(\sum_{i=1}^n \pi_i(a)v_i) = G(a') + K(a')(\sum_{i=1}^n \pi_i(a')v_i')$, i.e. the agent's expected utility is the same under both schemes, we must have*

$$(3.1) \qquad \sum_i [\pi_i(a') - \pi_i(a)](v_i' - v_i) \geq 0.$$

PROOF: From (2.2) and the assumption that the agent's expected utility is the same, we have

$$G(a') + K(a')\left( \sum_{i=1}^n \pi_i(a')v_i \right) \leq G(a) + K(a)\left( \sum_{i=1}^n \pi_i(a)v_i \right)$$

$$= G(a') + K(a')\left( \sum_{i=1}^n \pi_i(a')v_i' \right),$$

$$G(a) + K(a)\left( \sum_{i=1}^n \pi_i(a)v_i' \right) \leq G(a') + K(a')\left( \sum_{i=1}^n \pi_i(a')v_i' \right)$$

$$= G(a) + K(a)\left( \sum_{i=1}^n \pi_i(a)v_i \right).$$

It follows from the first of these that $\sum_{i=1}^{n} \pi_i(a')(v_i' - v_i) \geq 0$ and from the second that $\sum_{i=1}^{n} \pi_i(a)(v_i - v_i') \geq 0$ (since $K(a) > 0$ by Assumption A1(3)). Adding yields (3.1).                                                                              $Q.E.D.$

We now use Lemma 1 to show that an optimal incentive scheme will have the property that the principal's and agent's returns are positive related over some range of output levels; i.e., it is not optimal to have, for all output levels $q_i$, $q_j: I_i > I_j \Rightarrow q_i - I_i < q_j - I_j$. The proof proceeds by showing that, if the principal's and agent's payments are negatively related, then a twist in the incentive schedule which raises the agent's payment in high return states for the principal and lowers it in low return states for the principal can make the principal better off. The reason is that such a twist will be good for incentives since it gets the agent to put more probability weight on states yielding the principal a high return, and it is also good for risk-sharing since it raises the agent's return in low return states for the agent and lowers the agent's return in high return states for the agent. Since the incentive and risk-sharing effects reinforce each other, the principal is made better off.

In order to bring about both the incentive and risk-sharing effects, the twist in the incentive scheme must be chosen carefully. It is for this reason that the proof of the next proposition may seem rather complicated at first sight.

PROPOSITION 4: *Assume A1–A3 and V strictly concave. Let $(I_1, \ldots, I_n)$ be a second-best optimal incentive scheme. Then the following cannot be true: $I_i > I_j \Rightarrow q_i - I_i \leq q_j - I_j$ for all $1 \leq i, j \leq n$ and for some i, j, $I_i > I_j$ and $q_i - I_i < q_j - I_j$.*

PROOF: Suppose that

(3.2)        $I_i > I_j \Rightarrow q_i - I_i \leq q_j - I_j$

for all $1 \leq i, j \leq n$ and for some i, j, $I_i > I_j$ and $q_i - I_i < q_j - I_j$.

Let $(I_1', \ldots, I_n')$ be a new incentive scheme satisfying

(3.3)        $v_i' + \lambda h(v_i') = v_i + \lambda q_i - \mu$      for all $i$

where $v_i = V(I_i)$, $v_i' = V(I_i')$, $\lambda > 0$, and $\mu$ is such that

(3.4)        $\lambda \max_i (q_i - h(v_i)) \geq \mu \geq \lambda \min_i (q_i - h(v_i))$.

If $\lambda = \mu = 0$, then $v_i' = v_i$ solves (3.3). The implicit function theorem therefore implies that (3.3) has a solution as long as $\lambda$, $\mu$ are small. (Even if $h$ is not differentiable it has left and right hand derivatives.)

It follows from (3.2) and (3.4) that the change to the new incentive scheme has the effect of increasing the lowest $I_i$'s and decreasing the highest ones. For each $\lambda$ pick $\mu$ so that $G(a') + K(a')(\sum_{i=1}^{n} \pi_i(a')v_i') = \max_{a \in A}[G(a) + K(a)(\sum_{i=1}^{n} \pi_i(a) v_i')] = \max_{a \in A}[G(a) + K(a)(\sum_{i=1}^{n} \pi_i(a)v_i)]$. This ensures that the agent's expected

utility remains the same. We now show that the principal's expected profit is higher under the new incentive scheme than under the old, which contradicts the optimality of $(I_1, \ldots, I_n)$.

Substituting (3.1) of Lemma 1 into (3.3) yields:

$$\sum_i \pi_i(a')(q_i - h(v_i')) \geq \sum_i \pi_i(a)(q_i - h(v_i')).$$

If we can show that $\sum \pi_i(a)h(v_i') < \sum \pi_i(a)h(v_i)$, it will follow that

$$\sum \pi_i(a')(q_i - h(v_i')) > \sum \pi_i(a)(q_i - h(v_i)),$$

i.e., the principal is better off.

To see that $\sum \pi_i(a)h(v_i') < \sum \pi_i(a)h(v_i)$, note that

$$\sum \pi_i(a)(h(v_i) - h(v_i')) \geq \sum \pi_i(a)h'(v_i')(v_i - v_i')$$

by the convexity of $h$ (here $h'$ is the right-hand derivative if $h$ is not differentiable). It suffices therefore to show that the latter expression is positive. By (3.3),

$$\sum \pi_i(a)h'(v_i')(v_i - v_i') = \sum \pi_i(a)h'(v_i')(\lambda h(v_i') - \lambda q_i + \mu).$$

Suppose that this is nonpositive for small $\lambda$. Divide by $\lambda$ and let $\lambda \to 0$. Assuming without loss of generality $\mu/\lambda$ converges to $\hat{\mu}$ (we allow $\hat{\mu}$ infinite) and that $h'(v_i')$ converges to $\hat{h}_i'$, and using the fact that $v_i' \to v_i$, we get

$$(3.5) \qquad \sum \pi_i(a)\hat{h}_i'(h(v_i) - q_i + \hat{\mu}) \leq 0.$$

However, from the fact that $h'(v_i')$ is nondecreasing in $v_i'$ and $v_i' \to v_i$, $h'(v_i') \to \hat{h}_i'$, it follows that $v_i > v_j \Rightarrow \hat{h}_i' \geq \hat{h}_j'$. Hence by (3.2) $\hat{h}_i'$ and $(h(v_i) - q_i)$ are similarly ordered in the sense of Hardy, Littlewood, and Polya [5]; i.e., as one moves up so does the other. Therefore, by Hardy, Littlewood, and Polya [5, p. 43], $\hat{h}_i'$ and $(h(v_i) - q_i)$ are positively correlated, i.e.,

$$(3.6) \qquad \sum \pi_i(a)\hat{h}_i'(h(v_i) - q_i + \hat{\mu}) > \left(\sum \pi_i(a)\hat{h}_i'\right)\left(\sum \pi_i(a)(h(v_i) - q_i + \hat{\mu})\right)$$

$$\geq 0,$$

where the last inequality follows from the fact that (1) $h' \geq 0$; (2) $G(a) + K(a)$ $(\sum \pi_i(a)v_i') \leq G(a') + K(a')(\sum \pi_i(a')v_i') = G(a) + K(a)(\sum \pi_i(a)v_i)$ (since the agent's expected utility stays constant), which implies that

$$\lim_{\lambda \to 0}(1/\lambda)\sum \pi_i(a)(v_i - v_i') \geq 0.$$

(3.6) contradicts (3.5).

This proves that $\sum \pi_i(a)h(v_i') < \sum \pi_i(a)h(v_i)$, which establishes that the principal's expected profit is higher under $(I_1', \ldots, I_n')$. Contradiction.      $Q.E.D.$

REMARK 5: Another way of expressing Proposition 4 is that there is no

permutation $i_1, \ldots, i_n$ of the integers $1, \ldots, n$ such that $I_{i_k}$ is nondecreasing in $k$, and $(q_{i_k} - I_{i_k})$ is nonincreasing in $k$, with $I_{i_k} < I_{i_{k+1}}$, $(q_{i_k} - I_{i_k}) > (q_{i_{k+1}} - I_{i_{k+1}})$ for some $k$. Note that there is an interesting contrast between Proposition 4 and results found in the literature on optimal risk sharing in the absence of moral hazard. In this literature (see Borch [4]), it is shown that (if the individuals are risk averse) it is optimal for the individuals' returns to be positively related over the *whole* range of outcomes, whereas here we are only able to show that this is true over some range of outcomes.

Proposition 4 may be used to establish the following result about the monotonicity of the optimal incentive scheme.

PROPOSITION 5: *Assume A1–A3 and V strictly concave. Let $(I_1, \ldots, I_n)$ be a second-best optimal incentive scheme. Then* (1) *there exists $1 \le i \le n - 1$ such that $I_i \le I_{i+1}$, with strict inequality unless $I_1 = I_2 = \cdots = I_n$;* (2) *there exists $1 \le j \le n - 1$ such that $q_j - I_j < q_{j+1} - I_{j+1}$.*

PROOF: (1) follows directly from Proposition 4. So does (2) once we rule out the case $q_1 - I_1 = q_2 - I_2 = \cdots = q_n - I_n$. We do this by a similar argument to that used in Proposition 4. Suppose that $I$ is an optimal incentive scheme satisfying

$$(3.7) \qquad q_1 - I_1 = q_2 - I_2 = \cdots = q_n - I_n = k.$$

Then $I_1 < I_2 < \cdots < I_n$. Consider the new incentive scheme $I' = (I_1 + \epsilon, I_2 + \epsilon, \ldots, I_{n-1} + \epsilon, I_n - \mu\epsilon)$ where $\epsilon > 0$ and $\mu$ is chosen so that $\max_{a \in A}[G(a) + K(a)(\sum \pi_i(a)V(I_i'))] = \max_{A \in A}[G(a) + K(a)(\sum \pi_i(a)V(I_i))]$, i.e. the agent's expected utility is kept constant. We show that the principal's expected profit is higher under $I'$ than under $I$ for small $\epsilon$. Suppose not. Then

$$\sum \pi_i(a')(q_i - I_i') \le \sum \pi_i(a)(q_i - I_i) = k,$$

where $a'$ (resp. $a$) is optimal for the agent under $I'$ (resp. $I$). Substituting for $I'$ yields

$$-(1 - \pi_n(a'))\epsilon + \pi_n(a')\mu\epsilon \le 0.$$

Take limits as $\epsilon \to 0$. Without loss of generality $a' \to \hat{a}$. Hence we have

$$(3.8) \qquad -(1 - \pi_n(\hat{a})) + \pi_n(\hat{a})\mu \le 0.$$

Now since $a'$ is an optimal action for the agent under $I'$, it follows by uppersemicontinuity that $\hat{a}$ is optimal under $I$. Hence we have

$$G(\hat{a}) + K(\hat{a})\left(\sum \pi_i(\hat{a})V(I_i')\right) \le G(a') + K(a')\left(\sum \pi_i(a')V(I_i')\right)$$

$$= G(\hat{a}) + K(\hat{a})\left(\sum \pi_i(\hat{a})V(I_i)\right).$$

Hence $\sum \pi_i(\hat{a})(V(I_i) - V(I_i')) \ge 0$. Using the concavity of $V$ and taking limits as

$\epsilon \to 0$, we get

$$\sum_{i=1}^{n-1} \pi_i(\hat{a}) V'(I_i) - \pi_n(\hat{a}) V'(I_n) \mu \leq 0.$$

But since $V'(I_i)$ is decreasing in $i$, this contradicts (3.8). (If $V$ is not differentiable, $V'$ denotes the right-hand derivative.)

This proves that the principal does better under $I'$ than under $I$. Hence we have ruled out the case $q_1 - I_1 = \cdots = q_n - I_n$. This establishes Proposition 5.

$$Q.E.D.$$

Proposition 5 says that it is not optimal for the agent's marginal reward as a function of income to be negative everywhere or to be greater than or equal to one everywhere.[8] However, the proposition does allow for the possibility that either of these conditions can hold over some interval. To see when this may occur, it is useful to consider in more detail the case where $A$ is a finite set. When $A$ is finite, we can use the Kuhn–Tucker conditions for problem (2.2) to characterize the optimum. If Assumption A3 holds and $h$ is differentiable, these yield:

$$(3.9) \qquad h'(v_i) = \left[ \lambda + \sum_{\substack{a_j \in A \\ a_j \neq a^*}} \mu_j \right] K(a^*) - \sum_{\substack{a_j \in A \\ a_j \neq a^*}} \mu_j K(a_j) \left( \frac{\pi_i(a_j)}{\pi_i(a^*)} \right) \qquad \text{for all } i,$$

where $\lambda, (\mu_j)$ are nonnegative Lagrange multipliers and $\mu_j > 0$ only if the agent is indifferent between $a^*$ and $a_j$ at the optimum. The following proposition states that $\mu_j > 0$ for at least one action which is less costly than $a^*$. This implies that at an optimum the agent must be indifferent between at least two actions (unless $a^*$ is the least costly action, i.e. where there is no incentive problem).

PROPOSITION 6: *Assume* A1–A3 *and* $A$ *finite. Suppose that* (2.2) *has a solution for* $a^* \in A$. *Then if* $C_{FB}(a^*) > \min_{a \in A} C_{FB}(a)$, *this solution will have the property that* $G(a^*) + K(a^*)(\sum_{i=1}^n \pi_i(a^*)v_i) = G(a_j) + K(a_j)(\sum_{i=1}^n \pi_i(a_j)v_i)$ *for some* $a_j \in A$ *with* $C_{FB}(a_j) < C_{FB}(a^*)$. *Furthermore, if* $V$ *is strictly concave and differentiable, the Lagrange multiplier* $\mu_j$ *will be strictly positive for some* $a_j$ *with* $C_{FB}(a_j) < C_{FB}(a^*)$.

PROOF: Suppose that the agent strictly prefers $a^*$ to all actions less costly than $a^*$ at the solution. Then, since (2.2) is a convex problem, we can drop all the constraints in (2.2) which refer to less costly actions without affecting the

---

[8] Among other things, Proposition 5 shows that it is not optimal to have $q_1 - I_1 = q_2 - I_2 = \cdots = q_n - I_n$. This result has also been established by Shavell [20] under stronger assumptions.

solution. In other words, we can substitute $A' = \{a \in A \mid a$ is at least as costly as $a^*\}$ for $A$ in (2.2) and the solution will not change. But since $a^*$ is now the least costly action, we know from the proof of Proposition 1 that it is optimal to set $I_i = I_j$ for all $i, j$. However, $I_i = I_j$ is not optimal for the original problem since, under these conditions, the agent will pick an $a$ which minimizes $C_{FB}(a)$, and by assumption $C_{FB}(a^*) > \min_{a \in A} C_{FB}(a)$. Contradiction.

That $\mu_j > 0$ follows from the fact that if all the $\mu_j = 0$, then $h'(v_i)$ is the same for all $i$, which implies that $I_1 = \cdots = I_n$; however, this means that the agent will choose a cost-minimizing action, contradicting $C_{FB}(a^*) > \min_{a \in A} C_{FB}(a)$.

$$Q.E.D.$$

It should be noted that Proposition 6 depends strongly on the assumption that $A$ is finite.

The simplest case occurs when $\mu_j > 0$ for just one $a_j$ with $C_{FB}(a_j) < C_{FB}(a^*)$ (this will be true in particular if $A$ contains only two actions). In this case, we can rewrite (3.9) as

$$(3.10) \quad h'(v_i) = (\lambda + \mu_j)K(a^*) - \mu_j K(a_j)\frac{\pi_i(a_j)}{\pi_i(a^*)} \, .$$

We see that what determines $v_i$, and hence $I_i$, in this case is the relative likelihood that the outcome $q = q_i$ results from $a_j$ rather than from $a^*$. In particular, since $h$ convex $\Rightarrow h'$ nondecreasing in $v_i$, a sufficient condition for the optimal incentive scheme to be nondecreasing everywhere, i.e. $I_1 \leq I_2 \leq \cdots \leq I_n$, is that $\pi_i(a_j)/\pi_i(a^*)$ is nonincreasing in $i$, i.e. the relative likelihood that $a = a_j$ rather than $a = a^*$ produces the outcome $q = q_i$ is lower the better is the outcome $i$.

This observation has led some to suggest that the following is a sufficient conditon for the incentive scheme to be nondecreasing.

MONOTONE LIKELIHOOD RATIO CONDITION (MLRC): Assume A3. Then MLRC holds if, given $a, a' \in A$, $C_{FB}(a') \leq C_{FB}(a)$ implies that $\pi_i(a')/\pi_i(a)$ is nonincreasing in $i$.

It should be noted that the "first-order condition" approach described in the introduction, which is based on the assumption that the agent is indifferent between $a$ and $a + da$ at an optimum, does yield MLRC as a sufficient condition for monotonicity.[9] We now show, however, that, once we take into account the possibility that the agent may be indifferent between several actions at an

---

[9]See Mirrlees [11] or Holmstrom [7]. Milgrom [9] has shown that MLRC, as stated here, implies the differential version of the monotone likelihood condition which is to be found in Mirrlees [11] or Holmstrom [7].

optimum, i.e. $\mu_j > 0$ for more than one $a_j$, MLRC does not guarantee monotonicity.

EXAMPLE 1: $A = \{a_1, a_2, a_3\}$, $n = 3$. $\pi(a_1) = (\frac{2}{3}, \frac{1}{4}, \frac{1}{12})$, $\pi(a_2) = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$, $\pi(a_3) = (\frac{1}{12}, \frac{1}{4}, \frac{2}{3})$. Assume additive separability with $G(a_1) = 0$, $G(a_2) = -(\frac{1}{12}\sqrt{2} + \frac{1}{4}\sqrt{7/4})$, $G(a_3) = -\frac{7}{12}\sqrt{7/4}$, $V(I) = (3I)^{1/3}$ (i.e. $h(v) = \frac{1}{3}v^3), K(a) \equiv 1$ and $\overline{U} = \frac{1}{4}\sqrt{2} + \frac{1}{12}\sqrt{7/4}$. Note that MLRC is satisfied here.[10]

We compute $C(a_1), C(a_2), C(a_3)$. Obviously, $C(a_1) = C_{FB}(a_1) = \frac{1}{3}(\overline{U} - G(a_1))^3$ $= 0.033$. To compute $C(a_2)$, we use the first-order conditions (3.9). These are

$$v_1^2 = \lambda - \mu_1 + \tfrac{3}{4}\mu_2,$$

$$v_2^2 = \lambda + \tfrac{1}{4}\mu_1 + \tfrac{1}{4}\mu_2,$$

$$v_3^2 = \lambda + \tfrac{3}{4}\mu_1 - \mu_2,$$

plus the complementary slackness conditions. These equations are solved by setting $\lambda = \frac{5}{4}$, $\mu_1 = 2$, $\mu_2 = 1$. This yields $v_1 = 0$, $v_2 = \sqrt{2}$, $v_3 = \sqrt{7/4}$, and the agent is then indifferent between $a_1$, $a_2$, and $a_3$:

$$\tfrac{2}{3}v_1 + \tfrac{1}{4}v_2 + \tfrac{1}{12}v_3 + G(a_1) = \tfrac{1}{3}v_1 + \tfrac{1}{3}v_2 + \tfrac{1}{3}v_3 + G(a_2)$$

$$= \tfrac{1}{12}v_1 + \tfrac{1}{4}v_2 + \tfrac{2}{3}v_3 + G(a_3) = \overline{U}.$$

Since the first-order conditions are necessary and sufficient, we may conclude that $C(a_2) = \frac{1}{3}(\frac{1}{3}v_1^3 + \frac{1}{3}v_2^3 + \frac{1}{3}v_3^3) = 0.571$.

Note that the incentive scheme which implements $a_2$, $I_1 = 0$, $I_2 = \frac{1}{3}2^{3/2}$, $I_3 = \frac{1}{3}(\frac{7}{4})^{3/2}$, is not nondecreasing.

Observe that $C(a_3) \geq C_{FB}(a_3) = \frac{1}{3}(\overline{U} - G(a_3))^3 = 0.635 > C(a_2)$. Since $C(a_3) > C(a_2) > C(a_1)$, it is easy to show that we can find $q_1 < q_2 < q_3$ such that $B(a_2) - C(a_2) > \max[B(a_3) - C(a_3), B(a_1) - C(a_1)]$. But this means that it is optimal for the principal to get the agent to pick $a_2$. Hence the optimal incentive scheme is as described above. It is not nondecreasing despite the satisfaction of MLRC.

The reason that monotonicity breaks down in Example 1 is because, at the optimum, the agent is indifferent between $a_2$, the action to be implemented, $a_1$ a less costly action, and $a_3$ a more costly action. By MLRC $\pi_i(a_1)/\pi_i(a_2), \pi_i(a_2)/\pi_i(a_3)$ are decreasing in $i$. However, $\mu_1(\pi_i(a_1)/\pi_i(a_2)) + \mu_2(\pi_i(a_3)/\pi_i(a_2))$ need not be monotonic.

This observation suggests that one way to get monotonicity is to strengthen MLRC so that it holds for weighted combinations of actions as well as for the

---

[10]The function $V$ violates (2) of Assumption A1, but this is unimportant for the example.

basic actions themselves. In particular, suppose that

(3.11)    given any finite subset $\{a_1, \ldots, a_m\}$ of $A$, $a \in A$,

and nonnegative weights $w_1, \ldots, w_m$ summing to 1,

it is the case that $\left( \sum_{j=1}^{m} w_j \pi_i(a_j) / \pi_i(a) \right)$

is either nondecreasing in $i$ or nonincreasing in $i$.

Then, by the first-order conditions (3.9),

$$(3.12) \quad h'(v_i) = \left[ \lambda + \sum_{\substack{a_j \in A \\ a_j \neq a^*}} \mu_j \right] K(a^*) - \left[ \sum_{\substack{a_j \in A \\ a_j \neq a^*}} \mu_j K(a_j) \right] \left[ \sum_{\substack{a_j \in A \\ a_j \neq a^*}} w_j \left( \frac{\pi_i(a_j)}{\pi_i(a^*)} \right) \right],$$

where

$$w_j = \mu_j K(a_j) \Big/ \sum_{\substack{a_h \in A \\ a_h \neq a^*}} \mu_h K(a_h).$$

But, by (3.11), the right-hand side (RHS) of (3.12) is monotonic. Hence, the $v_i$'s are either monotonically nondecreasing or nonincreasing. By Proposition 5, however, they cannot be nonincreasing; hence they are nondecreasing.

Unfortunately, (3.11) turns out to be a very strong condition. In fact, it is equivalent to the following spanning condition.

SPANNING CONDITION (SC): There exists $\hat{\pi}, \hat{\pi}' \in S$ such that (1) for each $a \in A$, $\pi(a) = \lambda(a)\hat{\pi} + (1 - \lambda(a))\hat{\pi}'$ for some $0 \leq \lambda(a) \leq 1$; (2) $\hat{\pi}_i / \hat{\pi}_i'$ is nonincreasing in $i$.

That SC implies (3.11) is easy to see. We are grateful to Jim Mirrlees for pointing out and proving the converse.[11]

PROPOSITION 7: *Assume A1–A3, V strictly concave and differentiable. Suppose that SC holds. Then a second-best optimal incentive scheme satisfies* $I_1 \leq I_2 \leq \cdots \leq I_n$.

PROOF: If $A$ is finite, the argument following (3.12) establishes the result. To establish the result for the case $A$ infinite, let $\hat{a} \in A$ be a second-best optimal

---

[11] To prove the converse, define $a \lesssim a'$ if $\pi_i(a')/\pi_i(a)$ is nondecreasing in $i$. (3.11) implies that $\lesssim$ is a complete pre-ordering on $A$. Furthermore, $\lesssim$ is continuous. Since $A$ is compact, there exist $\underline{a}, \bar{a} \in A$ such that $\underline{a} \lesssim a \lesssim \bar{a}$ for all $a \in A$. Given $a \in A$, consider $\lambda(\pi_i(\bar{a})/\pi_i(a)) + (1 - \lambda)(\pi_i(\underline{a})/\pi_i(a))$. When $\lambda = 1$, this is nondecreasing in $i$, and when $\lambda = 0$, it is nonincreasing in $i$. Furthermore, (3.11) implies that it is monotonic in $i$ for all $0 < \lambda < 1$. It follows by continuity that it is independent of $i$ for some $0 < \lambda < 1$.

action and let $I$ be the second-best optimal incentive scheme which implements it. By Remark 2 of Section 2, $I$ is unique. Let $A_r$ be a finite subset of $A$ containing $\hat{a}$ such that the Euclidean distance between $A_r$ and $A$ is less than $(1/r)$. Let $I_r$ be the second-best optimal incentive scheme which implements $\hat{a}$ when the agent is restricted to choosing from $A_r$. From Proposition 7 for the finite $A$ case, we know that $I_r$ is nondecreasing. Take limits as $r \to \infty$. It is straightforward to show that $I_r \to I$. It follows that $I$ is nondecreasing.     Q.E.D.

An alternative sufficient condition for monotonicity may be found in the work of Mirrlees [12], who establishes a similar result to Proposition 8 below. For each $a \in A$, let $F(a) = (\pi_1(a), \pi_1(a) + \pi_2(a), \ldots, \pi_1(a) + \cdots + \pi_n(a))$. In the follow-ing proposition, the notation $F(a) \geq F'(a)$ is used to mean $F_i(a) \geq F_i'(a)$ for all $i = 1, \ldots, n$.

CONCAVITY OF DISTRIBUTION FUNCTION CONDITION (CDFC): CDFC holds if $a, a', a'' \in A$, and

$$\left( \frac{\bar{U} - G(a)}{K(a)} \right) = \lambda \left( \frac{\bar{U} - G(a')}{K(a')} \right) + (1 - \lambda) \left( \frac{\bar{U} - G(a'')}{K(a'')} \right),$$

$$0 \leq \lambda \leq 1,$$

imply that $F(a) \leq \lambda F(a') + (1 - \lambda) F(a'')$.

PROPOSITION 8: *Assume A1–A3, $V$ strictly concave and differentiable. Assume also that $U$ is additively or multiplicatively separable, i.e., either $G(a) \equiv 0$ or $K(a) \equiv constant$. Suppose that MLRC and CDFC hold. Then a second-best optimal incentive scheme $(I_1, \ldots, I_n)$ satisfies $I_1 \leq I_2 \leq \cdots \leq I_n$.*

PROOF: Assume first that $A$ is finite. Let $a^*$ maximize $B(a) - C(a)$. Let $A' = \{ a \in A \mid C_{FB}(a) \leq C_{FB}(a^*) \}$. Consider the cost minimizing way of getting the agent to pick $a^*$ given that he can choose only from $A'$. It is clear from (3.9) that, since $\pi_i(a_j)/\pi_i(a^*)$ is nonincreasing in $i$ by MLRC, the incentive scheme $(I_1, \ldots, I_n)$ is nondecreasing. We will be home if we can show that $(I_1, \ldots, I_n)$ is optimal when $A'$ is replaced by $A$. Since adding actions cannot reduce the cost of implementing $a^*$, all we have to do is to show that $(I_1, \ldots, I_n)$ continues to implement $a^*$, i.e. there does not exist $a''$, $C_{FB}(a'') > C_{FB}(a^*)$, such that

$$(3.13) \quad G(a'') + K(a'')\left( \sum \pi_i(a'')v_i \right) > G(a^*) + K(a^*)\left( \sum \pi_i(a^*)v_i \right).$$

However, we know from Propositions 2 and 6 that

$$(3.14) \quad G(a^*) + K(a^*)\left( \sum \pi_i(a^*)v_i \right) = G(a') + K(a')\left( \sum \pi_i(a')v_i \right) = \bar{U}$$

for some $a'$ with $C_{FB}(a') < C_{FB}(a^*)$. Writing

$$\frac{\bar{U} - G(a^*)}{K(a^*)} = \lambda \left( \frac{\bar{U} - G(a'')}{K(a'')} \right) + (1 - \lambda) \left( \frac{\bar{U} - G(a')}{K(a')} \right)$$

and using CDFC and the fact that $v_1 \le v_2 \le \cdots \le v_n$, we get

$$\sum \pi_i(a^*)v_i - \left( \frac{\overline{U} - G(a^*)}{K(a^*)} \right)$$

$$\ge \lambda \sum \pi_i(a'')v_i + (1-\lambda)\left( \sum \pi_i(a')v_i \right) - \left( \frac{\overline{U} - G(a^*)}{K(a^*)} \right)$$

$$= \lambda \left[ \sum \pi_i(a'')v_i - \left( \frac{\overline{U} - G(a'')}{K(a'')} \right) \right]$$

$$+ (1-\lambda)\left[ \sum \pi_i(a')v_i - \left( \frac{\overline{U} - G(a')}{K(a')} \right) \right].$$

But this contradicts (3.13) and (3.14).

To prove the result for $A$ finite, one again proceeds by way of finite approximation. *Q.E.D.*

To understand CDFC, consider, for each $a \in A$, $V(C_{FB}(a)) = ((\overline{U} - G(a)) / K(a))$. In utility terms $V(C_{FB}(a))$ is a measure of the first-best cost of getting the agent to pick $a$. CDFC says that if $a$ is a convex combination of $a'$ and $a''$ in terms of this measure of cost then the distribution function of outcomes corresponding to $a$ dominates in the sense of first degree stochastic dominance the corresponding convex combination of the distribution functions corresponding to $a'$ and $a''$. It is worth noting that under the assumption of additive or multiplicative separability in Proposition 8, the $\lambda$ in the CDFC definition is independent of $\overline{U}$.

So far we have considered only the monotonicity of the optimal incentive scheme. One would also like to know when the optimal incentive scheme is *progressive*, i.e. $(I_{i+1} - I_i)/(q_{i+1} - q_i)$ is nonincreasing in $i$, or *regressive*, i.e. $(I_{i+1} - I_i)/(q_{i+1} - q_i)$ is nondecreasing in $i$. To get results about this, one needs considerably stronger assumptions, as the following proposition indicates.

PROPOSITION 9: *Assume A1–A3, $V$ strictly concave and differentiable. Assume also that $U$ is additively or multiplicatively separable, i.e., either $G(a) \equiv 0$ or $K(a) \equiv$ constant. Suppose that MLRC and CDFC hold and that $(q_{i+1} - q_i)$ is independent of $i$, $1 \le i \le n-1$. Then a second-best optimal incentive scheme will be regressive (resp. progressive) if*

(3.15)  $(1/V'(I))$ *is concave (resp. convex) in $I$ and $a, a' \in A$,*

$C_{FB}(a') < C_{FB}(a)$, *implies that* $(\pi_{i+1}(a')/\pi_{i+1}(a)) - (\pi_i(a')/\pi_i(a))$

*is nonincreasing (resp. nondecreasing) in $i$.*

PROOF: Assume first that $A$ is finite. Let $a^*$ be a second-best optimal action. Let $a'$ maximize $C_{FB}(a)$ subject to $C_{FB}(a) < C_{FB}(a^*)$, i.e. $a'$ is the next most costly action after $a^*$. Consider the cost minimizing way of implementing $a^*$ given that $a'$ is the only other action that the agent can choose. Using the same concavity argument as in the proof of Proposition 8, we can show that the resulting incentive scheme $(I_1, \ldots, I_n)$ also implements $a^*$ when the agent can choose from all of $A$. Hence $(I_1, \ldots, I_n)$ is an optimal incentive scheme.

By (3.10),

$$\frac{1}{V'(I_i)} = h'(v_i) = (\lambda + \mu)K(a^*) - \mu K(a') \frac{\pi_i(a')}{\pi_i(a^*)}$$

and so

$$\frac{1}{V'(I_{i+1})} - \frac{1}{V'(I_i)} = -\mu K(a') \left( \frac{\pi_{i+1}(a')}{\pi_{i+1}(a^*)} - \frac{\pi_i(a')}{\pi_i(a^*)} \right).$$

(3.15) now follows immediately. To prove the result for the $A$ infinite case, one again proceeds by way of a finite approximation.                    Q.E.D.

Note that $1/V'$ is linear if $V = \log I$; is concave if $V = -e^{-\alpha I}$, $\alpha > 0$, or $V = I^\alpha$, $0 < \alpha < 1$; is convex if $V = -I^{-\alpha}$, $\alpha > 1$.

It should also be noted that Mirrlees [12] has shown that if CDFC holds, the "first-order condition" approach referred to in the introduction is valid. Thus Propositions 8 and 9 can also be proved by appealing to the characterization of an optimal incentive scheme to be found in much of the literature (see, e.g., Holmstrom [7] and Mirrlees [11]).

Let us summarize the results of this section. We have shown that an optimal incentive scheme will not be declining everywhere, but that only under quite strong assumptions (SC or MLRC plus concavity) will it be nondecreasing everywhere. We have also shown that it is not optimal for the agent's marginal remuneration for an extra pound of profit to exceed one everywhere, although it may exceed one sometimes. Finally, we have obtained sufficient conditions for the incentive scheme to be progressive or regressive.

The conclusion that only under strong assumptions will the optimal incentive scheme be monotonic may seem disappointing at first sight. One feels that monotonicity is a minimal requirement. This may not be the right reaction, however. There are many interesting situations where it is clear that the optimal scheme will not be monotonic. We have described one example in the introduction. Another example is the following. Suppose that actions are two dimensional, with one dimension referring to how hard the agent works and the other dimension to how cautious he is—greater caution might lead to a lower variance of profit but also to a lower mean. The optimal action for the principal might involve the agent working fairly hard and also not being too cautious. The best

way to implement this may be to pay the agent high amounts for both very good outcomes (to encourage high effort) and very bad outcomes (to discourage excessive caution). This example seems far from pathological. In fact, one might argue that a number of real world incentive schemes operate in this way. In view of examples like this, the difficulty of finding general conditions guaranteeing monotonicity may become less surprising.[12]

In the next section, we show that considerably stronger results than those of this section can be proved for the case $n = 2$. We also provide a simple algorithm for computing optimal incentive schemes when $n = 2$.

### 4. THE CASE OF TWO OUTCOMES

When $n = 2$, we will refer to $q_1$ as the "bad" outcome and $q_2 > q_1$ as the "good" outcome. In this case, the agent's incentive scheme can be represented simply by a fixed payment $w$ and a share of profits, $s$, where $w + sq_1 = I_1$, $w + sq_2 = I_2$, i.e., $s = (I_2 - I_1)/(q_2 - q_1)$. Proposition 5 of the last section shows that it is not optimal for $I_i$ to be everywhere declining in $q_i$. When $n = 2$, this means that $s \geq 0$.[13] Similarly the proposition implies that $s < 1$ when $n = 2$. This has a number of interesting implications.

DEFINITION: Let $n = 2$. We say that $a \in A$ is *efficient* if there does not exist $a' \in A$ satisfying $C_{FB}(a') \leq C_{FB}(a)$ and $\pi_2(a') \geq \pi_2(a)$, with at least one strict inequality.

In other words, an action is efficient if the probability of a good outcome can only be increased by incurring greater cost.

PROPOSITION 10: *Assume A1–A3 and V strictly concave. Let $n = 2$. Then every second-best optimal action is efficient.*

PROOF: Let $a$ be a second-best optimal action. Then $a$ maximizes $G(a) + K(a)$ $[\pi_1(a)v_1 + \pi_1(a)v_2]$. Suppose $C_{FB}(a') \leq C_{FB}(a)$ and $\pi_2(a') \geq \pi_2(a)$, with at least one strict inequality. Then, by the definition of $C_{FB}$,

$$G(a) + K(a)V(C_{FB}(a)) = \overline{U} = G(a') + K(a')V(C_{FB}(a'))$$

$$\leq G(a') + K(a')V(C_{FB}(a)),$$

---

[12] There are some cases where monotonicity may be a *constraint* on the optimal incentive scheme. An example is where the agent can always make a better outcome look like a worse outcome by reducing the firm's profits after the outcome has occurred. This case can be analyzed by adding the (linear) constraints $v_1 \leq v_2 \leq \cdots \leq v_n$ to the problem (2.2).

[13] Shavell [19] also proves that $s \geq 0$ when $n = 2$, but under stronger assumptions.

since $C_{FB}(a') \leq C_{FB}(a)$. Hence, by Assumption A1(4), $G(a) + K(a)v \leq G(a') + K(a')v$ for all $v \in \mathfrak{U}$. Therefore using the fact that $v_1 \leq v_2$ since $s \geq 0$, and the fact that $\pi_2(a') \geq \pi_2(a)$, we have

$$G(a) + K(a)\big[\pi_1(a)v_1 + \pi_2(a)v_2\big]$$

$$\leq G(a') + K(a')\big[\pi_1(a)v_1 + \pi_2(a)v_2\big]$$

$$\leq G(a') + K(a')\big[\pi_1(a')v_1 + \pi_2(a')v_2\big]$$

with at least one strict inequality unless $C_{FB}(a) = C_{FB}(a')$ and $v_1 = v_2$. This contradicts the optimality of $a$ unless $C_{FB}(a) = C_{FB}(a')$ and $v_1 = v_2$. However, in this case, the agent is indifferent between $a$ and $a'$, while the principal prefers $a'$, again contradicting the optimality of $a$.                    $Q.E.D.$

We may use Proposition 10 to prove that when $n = 2$ it will never pay the principal to offer the agent an expected utility in excess of $\overline{U}$ (recall that when $n > 2$ this is only generally true when $U(a, I)$ is additively or multiplicatively separable—see Proposition 2).

PROPOSITION 11: *Assume A1–A3 and V strictly concave. Let $n = 2$. Let $\hat{a}$ be a second-best optimal action and $\hat{I}$ a second-best optimal incentive scheme which implements $\hat{a}$. Then $\sum_{i=1}^{n} \pi_i(\hat{a}) U(\hat{a}, \hat{I}_i) = \overline{U}$.*

PROOF: Suppose not, i.e., $\sum_{i=1}^{n} \pi_i(\hat{a}) U(\hat{a}, \hat{I}_i) > \overline{U}$. Consider a new incentive scheme $(I_1, I_2) = (\hat{I}_1 - \epsilon, \hat{I}_2)$ where $\epsilon > 0$ is small. Let $a$ be an optimal action for the agent under the new scheme, i.e., $a$ maximizes $G(a) + K(a)[\pi_1(a)V(\hat{I}_1 - \epsilon) + \pi_2(a)V(\hat{I}_2)]$. Then,

$$\pi_1(a)(q_1 - I_1 + \epsilon) + \pi_2(a)(q_2 - I_2) > \pi_1(a)(q_1 - I_1) + \pi_2(a)(q_2 - I_2)$$

$$\geq \pi_1(\hat{a})(q_1 - \hat{I}_1) + \pi_2(\hat{a})(q_2 - \hat{I}_2)$$

as long as $\pi_2(\hat{a}) \leq \pi_2(a)$ (since $0 \leq s < 1$). Thus, if we can show that $\pi_2(\hat{a}) \leq \pi_2(a)$, we will have contradicted the optimality of $(\hat{I}_1, \hat{I}_2)$, since the principal's profits will be higher under $(I_1, I_2)$ than under $(\hat{I}_1, \hat{I}_2)$.

Suppose $\pi_2(\hat{a}) > \pi_2(a)$. Now the same argument as in Proposition 10 shows that $a$ is efficient. Thus we must have $C_{FB}(\hat{a}) > C_{FB}(a)$. Hence $G(a) + K(a)V(C_{FB}(a)) = \overline{U} = G(\hat{a}) + K(\hat{a})V(C_{FB}(\hat{a})) > G(\hat{a}) + K(\hat{a})V(C_{FB}(a))$, and so, by Assumption A1(4),

(4.1)        $G(a) + K(a)v > G(\hat{a}) + K(\hat{a})v$

for all $v \in \mathfrak{U} \equiv \{V(I) \mid I \in \mathcal{I}\}$. Since $\mathfrak{U}$ contains arbitrarily large negative num-

bers, we may conclude from (4.1) that $K(a) \leq K(\hat{a})$. Now by revealed preference,

(4.2) $\quad G(a) + K(a)\left[\pi_1(a)V(\hat{I}_1) + \pi_2(a)V(\hat{I}_2)\right]$

$$\leq G(\hat{a}) + K(\hat{a})\left[\pi_1(\hat{a})V(\hat{I}_1) + \pi_2(\hat{a})V(\hat{I}_2)\right],$$

(4.3) $\quad G(a) + K(a)\left[\pi_1(a)V(\hat{I}_1 - \epsilon) + \pi_2(a)V(\hat{I}_2)\right]$

$$\geq G(\hat{a}) + K(\hat{a})\left[\pi_1(\hat{a})V(\hat{I}_1 - \epsilon) + \pi_2(\hat{a})V(\hat{I}_2)\right].$$

Subtracting (4.3) from (4.2) yields $K(a)\pi_1(a) \leq K(\hat{a})\pi_1(a)$. Hence, since $\pi_2(\hat{a}) > \pi_2(a)$ by assumption, $K(a) < K(\hat{a})$. However, rewriting (4.2), we obtain

$$G(a) + K(a)\bar{v} + K(a)\left[\pi_1(a)(V(\hat{I}_1) - \bar{v}) + \pi_2(a)(V(\hat{I}_2) - \bar{v})\right]$$

$$\leq G(\hat{a}) + K(\hat{a})\bar{v} + K(\hat{a})\left[\pi_1(\hat{a})(V(\hat{I}_1) - \bar{v}) + \pi_2(\hat{a})(V(\hat{I}_2) - \bar{v})\right]$$

where $\bar{v} = \sup \mathfrak{V}$. (Note that $\bar{v} < \infty$, for $\bar{v} = \infty$ and $K(a) < K(\hat{a})$ violate (4.1).) Setting $v = \bar{v}$ in (4.1), we may conclude that

$$K(a)\pi_1(a)(V(\hat{I}_1) - \bar{v}) + K(a)\pi_2(a)(V(\hat{I}_2) - \bar{v})$$

$$\leq K(\hat{a})\pi_1(\hat{a})(V(\hat{I}_1) - \bar{v}) + K(\hat{a})\pi_2(\hat{a})(V(\hat{I}_2) - \bar{v}).$$

But this is impossible since $K(a)\pi_1(a) \leq K(\hat{a})\pi_1(\hat{a})$, $K(a) < K(\hat{a})$, $\pi_2(a) < \pi_2(\hat{a})$, $V(\hat{I}_1) - \bar{v} < 0$, $V(\hat{I}_2) - \bar{v} < 0$. We have thus shown that $\pi_2(a) \geq \pi_2(\hat{a})$, which contradicts the optimality of $(\hat{I}_1, \hat{I}_2)$. $\hspace{2em} Q.E.D.$

Proposition 11 tells us that the agent's fixed payment $w$ is determined once $s$ is. In particular, $w$ will be the unique solution of

$$\max_{a \in A}\left[G(a) + K(a)(\pi_1(a)V(w + sq_1) + \pi_2(a)V(w + sq_2))\right] = \overline{U}.$$

We have shown that one implication of Proposition 5 for the case $n = 2$ is that every second-best optimal action is efficient. We consider now a second implication. Suppose that we start off in the situation where the agent has access to a set of actions $A$, and now some additional actions become available, so that the new action set is $A' \supset A$. Then, if the new actions are all higher cost actions for the agent than those in $A$—in the sense that their $C_{FB}$'s are higher—the principal cannot be made worse off by such a change.

PROPOSITION 12: *Assume* A1 *and* A2. *Let* $n = 2$. *Suppose that* $A' \supset A$ *and that* $a \in A$, $a' \in A' \backslash A \Rightarrow C_{FB}(a') \geq C_{FB}(a)$. *Assume that* A3 *holds for both* $A$ *and* $A'$. *Then* $\max_{a \in A}[B(a) - C'(a)] \geq \max_{a \in A}[B(a) - C(a)]$, *where* $C'$ *is the second-best cost function under* $A'$.

PROOF: Suppose $(I_1, I_2)$ is an optimal second-best incentive scheme when the action set is $A$. Let the principal keep this incentive scheme when the new actions

$A' \backslash A$ are added. The only way that the principal can be made worse off is if the agent now switches from $a \in A$ to $a' \in A' \backslash A$. But $a'$ must then provide higher utility for the agent, i.e., $G(a') + K(a')[\pi_1(a')v_1 + \pi_2(a')v_2] > G(a) + K(a)[\pi_1(a)v_1 + \pi_2(a)v_2]$. Since $C_{FB}(a') \geq C_{FB}(a)$, however, $G(a') + K(a')v \leq G(a) + K(a)v$ for all $v \in \mathcal{U}$ (by Assumption A1(4)). Hence $\pi_1(a')v_1 + \pi_2(a')v_2 > \pi(a)v_1 + \pi_2(a)v_2$, which implies, since $v_2 \geq v_1$ by Proposition 5, that $\pi_2(a') > \pi_2(a)$. But it follows that the principal's expected profits $\pi_1(q_1 - I_1) + \pi_2(q_2 - I_2)$ rise when the agent moves from $a$ to $a'$ since, again by Proposition 5, $s < 1$, i.e. $q_2 - I_2 > q_1 - I_1$.    Q.E.D.

As a final implication of Proposition 5, when $n = 2$, consider a manager-entrepreneur who initially owns 100 per cent of a firm, i.e. $\tilde{w} = 0$, $\tilde{s} = 1$. In the absence of any risk-sharing possibilities the manager will choose $a$ to maximize $\pi_1(a)U(a, q_1) + \pi_2(a)U(a, q_2)$. Let $\tilde{a}$ be a solution to this. Clearly $\tilde{a}$ is efficient. Now suppose a risk neutral principal appears with whom the manager can share risks. We know from Proposition 5 that at the new optimum $s < 1 = \tilde{s}$. Therefore, by Lemma 1 and Proposition 11, $\pi_2(a^*) \leq \pi_2(\tilde{a})$. In addition, $C_{FB}(a^*) \leq C_{FB}(\tilde{a})$ by Proposition 10. Thus, the existence of risk-sharing possibilities leads the agent to choose a less costly action with a lower probability of a good outcome.

We may use Propositions 10–12 to develop a method for computing a second-best optimal incentive scheme when $n = 2$. Consider the case where $A$ is finite. Recall that Proposition 6 states that, in this case, the agent will be indifferent between $a^*$ and some less costly action. This fact makes the computation of an optimal incentive scheme fairly straightforward. We know from Proposition 10 that it is never optimal to get the agent to choose an inefficient action. Hence we can assume without loss of generality that $C_{FB}(a_1) < C_{FB}(a_2) < \cdots < C_{FB}(a_m)$ and $\pi_2(a_1) < \pi_2(a_2) < \cdots < \pi_2(a_m)$. The computation of $C(a_1)$ is easy: by Remark 1 of Section 2 it is just $C_{FB}(a_1)$. To compute $C(a_k)$, $k > 1$, we use Propositions 6 and 11. For each action $a_j$, $j < k$, find $I_1, I_2$ so that the agent is indifferent between $a_k$ and $a_j$ and the agent's expected utility is $\overline{U}$. This means solving

$$
\begin{aligned}
G(a_k) + K(a_k)(\pi_1(a_k)v_1 + \pi_2(a_k)v_2) &= \overline{U}, \\
G(a_j) + K(a_j)(\pi_1(a_j)v_1 + \pi_2(a_j)v_2) &= \overline{U},
\end{aligned}
$$
(4.4)

which yields

$$
v_1 = \frac{\pi_2(a_j)\big((\overline{U} - G(a_k))/K(a_k)\big) - \pi_2(a_k)\big((\overline{U} - G(a_j))/K(a_j)\big)}{\pi_1(a_k) - \pi_1(a_j)},
$$
(4.5)
$$
v_2 = \frac{\pi_1(a_j)\big((\overline{U} - G(a_k))/K(a_k)\big) - \pi_1(a_k)\big((\overline{U} - G(a_j))/K(a_j)\big)}{\pi_2(a_k) - \pi_2(a_j)}.
$$

We then set $I_1 = h(v_1)$, $I_2 = h(v_2)$. Note that $v_1 < v_2$ in (4.5) so that $I_1 < I_2$.

Plot with vertical axis $v_2$ and horizontal axis $v_1$.

$\pi_1(a_{j2})v_1 + \pi_2(a_{j2})v_2 = (\overline{U} - G(a_{j2}))\,/\,K(a_{j2})$

$\pi_1(a_{j1})v_1 + \pi_2(a_{j1})v_2 = (\overline{U} - G(a_{j1}))\,/\,K(a_{j1})$

$\pi_1(a_k)v_1 + \pi_2(a_k)v_2 = (\overline{U} - G(a_k))\,/\,K(a_k)$

$45°$

FIGURE 2.

Doing this for each $j = 1, \ldots, k - 1$ yields $(k - 1)$ different $(v_1, v_2)$ (and $(I_1, I_2)$) pairs, each with $v_1 < v_2$. This is illustrated in Figure 2 for the case $k = 3$, where the $(v_1, v_2)$ pairs are at $A, B$. We know from Proposition 6 that one of these pairs is the solution to (2.2). In fact, the solution must occur at the $(v_1, v_2)$ pair with the *smallest* $v_1$ (and hence, by (4.4), with the largest $v_2$)—denote this pair by $(\hat{v}_1, \hat{v}_2)$. To see this, suppose that the agent is indifferent between $a_k$ and $a_j$ under $(\hat{v}_1, \hat{v}_2)$. Consider the expression

$$(4.6) \qquad \pi_1(a_k)v_1 + \pi_2(a_k)v_2 - \pi_1(a_j)v_1 - \pi_2(a_j)v_2$$

$$= (\pi_1(a_k) - \pi_1(a_j))v_1 + (\pi_2(a_k) - \pi_2(a_j))v_2 .$$

When $v_1 = \hat{v}_1$, $v_2 = \hat{v}_2$, this expression equals $[(\overline{U} - G(a_k))/K(a_k)] - [(\overline{U} - G(a_j))/K(a_j)]$. Suppose now that $v_1 > \hat{v}_1$, $v_2 < \hat{v}_2$. Then (4.6) falls since $\pi_1(a_k) < \pi_1(a_j)$. Hence the agent now prefers $a_j$ to $a_k$ and so $a_k$ is not implemented.

In Figure 2, the solution is at $A$. (The solution could not be at $B$ since it is clear from the diagram that, at $B$, $a_{j2}$ gives the agent an expected utility greater than $\overline{U}$, i.e. $a_k$ is not implemented at $B$.) Note that it is possible that the $(\hat{v}_1, \hat{v}_2)$ picked in this way does not lie in $\mathfrak{A} \times \mathfrak{A}$; i.e., $h(\hat{v}_1)$ or $h(\hat{v}_2)$ may be undefined. In this case, the constraint set of (2.2) is empty and so $C(a_k) = \infty$. If $(\hat{v}_1, \hat{v}_2) \in \mathfrak{A} \times \mathfrak{A}$, then the principal's minimum expected cost of getting the agent to pick $a_k$ is $C(a_k) = \pi_1(a_k)h(\hat{v}_1) + \pi_1(a_k)h(\hat{v}_2)$. The expected net benefits of implementing $a_k$ are $B(a_k) - C(a_k)$. This procedure must be undergone for each $a_k$,

$k = 1, \ldots, m$. Finally, the overall optimum is determined by selecting the $a$ which maximizes $B(a_k) - C(a_k)$.

REMARK 6: In computing the cost of implementing $a_k$, we have ignored actions which are more costly for the agent than $a_k$. This means that the cost function which we have computed is not the true cost function $C(a)$ but a modified cost function $\tilde{C}(a)$. Clearly, $\tilde{C}(a) \leq C(a)$ for each $a$ since more actions can only make implementation more difficult. On the other hand, Proposition 12 tells us that $\max_{a \in A}[B(a) - \tilde{C}(a)] \leq \max_{a \in A}[B(a) - C(a)]$. Combining these yields $\max_{a \in A}[B(a) - C(a)] = \max_{a \in A}[B(a) - \tilde{C}(a)]$, which means that we are justified in working with $\tilde{C}(a)$ instead of $C(a)$.

Another case where computation is quite simple is when $A$ is infinite and $\{C_{FB}(a) \mid a \in A\}$ is an interval $[\underline{c}, \bar{c}]$ of the real line. For reasons of space, we do not cover this case.

Unfortunately, the computational techniques presented above do not appear to generalize in a useful way to the case $n > 2$. In order to compute an optimum when $n > 2$, in the finite action case, it seems that we must, for each $a \in A$, solve the convex problem in (2.2) and then, by inspection, find the $a \in A$ which maximizes $B(a) - C(a)$. If $A$ is infinite, one takes a finite approximation. These steps can be carried out on a computer, although the amount of computer time involved when the number of elements of $A$ is large may be considerable.

One case where a considerable simplification can be achieved when $n > 2$ is where MLRC and CDFC hold. Then the solution to (2.2) has the property that (1) if $A$ is finite, the agent is indifferent only between $a^*$, the action the principal wants to implement, and $a'$, where $a'$ maximizes $C_{FB}(a)$ subject to $C_{FB}(a) < C_{FB}(a^*)$, i.e. $a'$ is the next most costly action after $a^*$ (see the proof of Proposition 8); (2) if $A$ is convex, then $a^*$ is the unique maximizer of $G(a) + K(a)(\sum \pi_i(a) V(I_i))$, and $[d(G(a^*) + K(a^*)(\sum \pi_i(a^*) V(I_i)))/da] = 0$ is a necessary and sufficient condition for the agent to pick $a^*$. In the latter case, Mirrlees [12] has shown that the first-order condition approach referred to in the introduction is valid.

One may ask also whether Propositions 10 and 12 hold in the case $n > 2$. The answer is no (but see Remark 7 below). Second-best optimal actions may be inefficient; i.e., there may exist lower cost actions which dominate the optimal action in the sense of first degree stochastic dominance.[14] Also the addition of actions costlier than the second-best optimal action may make the principal worse off (in Example 1, the principal's expected profits increase if action $a_3$

---

[14] Let $A = \{a_1, a_2, a_3\}$, $n = 3$. Assume $C_{FB}(a_1) < C_{FB}(a_2) < C_{FB}(a_3)$, and that $\pi(a_1) = (3/4, 1/8, 1/8)$, $\pi(a_2) = (1/3, 1/3, 1/3)$, $\pi(a_3) = (1/2, 1/2, 0)$ (Assumption A3 is violated, but this is unimportant.) Then $C(a_1) = C_{FB}(a_1)$ since $a_1$ is the least cost action, and $C(a_3) = C_{FB}(a_3)$ since $a_3$ can be implemented by setting $I_1 = I_2$, $I_3 = -\infty$. However, $C(a_2) > C_{FB}(a_2)$ and, in fact, if the agent is very risk averse, $C(a_2)$ will be so big that it is profitable for the principal to implement $a_3$ rather than $a_2$ (the effect of risk aversion on $C(a)$ is discussed in Section 5). This is in spite of the fact that $a_3$ is inefficient relative to $a_2$.

becomes unavailable to the agent). Finally, as Shavell [19] has noted the agent may choose a higher cost action when there are opportunities to share risks with a principal than in the absence of these opportunities.

REMARK 7: It is interesting to note that it is possible to extend all the results of the $n = 2$ case to the $n > 2$ case when the spanning condition (SC) holds. This is because when SC holds, both the principal and the agent are essentially choosing between lotteries of the probability vectors $\hat{\pi}$ and $\hat{\pi}'$.

In particular, let $I_1(v_1) = \min_{\{I_i\}} \sum_{i=1}^{n} \hat{\pi}_i I_i$ subject to $\sum_{i=1}^{n} \hat{\pi}_i V(I_i) \geq v_1; I_2(v_2)$ $= \min_{\{I_i\}} \sum_{i=1}^{n} \hat{\pi}'_i I_i$ subject to $\sum_{i=1}^{n} \hat{\pi}'_i V(I_i) \geq v_2$. Now consider the principal's minimum cost problem as: for each $a^*$, choose $v_1$ and $v_2$ to minimize $\lambda(a^*)$ $I_1(v_1) + (1 - \lambda(a^*))I_2(v_2)$ subject to (1) $G(a^*) + [\lambda(a^*)v_1 + (1 - \lambda(a^*))v_2]K(a^*)$ $\geq G(a) + [\lambda(a)v_1 + (1 - \lambda(a))v_2]K(a)$ for all $a \in A$; (2) $G(a^*) + [\lambda(a^*)v_1 + (1 - \lambda(a^*))v_2]K(a^*) \geq \overline{U}$. Then the principal's problem looks exactly the same as in the $n = 2$ case. Note that from stochastic dominance (i.e. part (2) of the SC condition) $\sum_{i=1}^{n} \pi_i q_i \leq \sum_{i=1}^{n} \pi'_i q_i$, so "state 2" is the good state. We are grateful to Bengt Holmstrom for alerting us to the fact that all of the results for the $n = 2$ case hold when $n > 2$ and the Spanning Condition is satisfied.

## 5. WHAT DETERMINES HOW SERIOUS THE INCENTIVE PROBLEM IS?

In previous sections, we have studied the properties of an optimal incentive scheme. We turn now to a consideration of the factors which determine the magnitude of $L$, the loss to the principal from being unable to observe the agent's action.

One feels intuitively that the worse is the quality of the information about the agent's action that the principal obtains from observing any outcome, the more serious will be the incentive problem. This idea can be formalized as follows. Suppose that we start with an incentive problem in which the agent's action set is $A$, his utility function is $U$, his reservation utility is $\overline{U}$, the probability function is $\pi$, and the vector of outputs is $q = (q_1, \ldots, q_n)$. We denote this incentive problem by $(A, U, \overline{U}, \pi, q)$. Consider the new incentive problem $(A, U, \overline{U}, \pi', q')$ where $\pi'(a) = R\pi(a)$ for all $a \in A$ and $R$ is an $(n \times n)$ stochastic matrix (here $\pi(a), \pi'(a)$ are $n$ dimensional column vectors and the columns of $R$ sum to one). Below we show that $C'(a) \geq C(a)$ for all $a \in A$, where unprimed variables refer to the original incentive problem and primed variables to the new incentive problem.

The transformation from $\pi(a)$ to $R\pi(a)$ corresponds to a decrease in informativeness in the sense of Blackwell (see, e.g., Blackwell and Girshick [3]).[15] That is, if we think of the actions $a \in A$ as being parameters with respect to which we

---

[15] The possibility of using Blackwell's notion of informativeness to characterize the seriousness of an incentive problem was suggested by Holmstrom [7].

have a prior probability distribution, then an experimenter who makes deductions about $a$ from observing $q_1, \ldots, q_n$ would prefer to face the function $\pi$ than the function $R\pi$.

PROPOSITION 13: *Consider the two incentive problems* $(A, U, \overline{U}, \pi, q)$, $(A, U, \overline{U}, \pi', q')$ *and assume that Assumptions* A1–A3 *hold for both. Suppose that* $\pi'(a)$ $= R\pi(a)$ *for all* $a \in A$, *where* $R$ *is an* $(n \times n)$ *stochastic matrix. Then* $C'(a)$ $\geq C(a)$ *for all* $a \in A$. *Furthermore, if* $V$ *is strictly concave and* $R \gg 0$,[16] *then* $C_{FB}(a^*) > \min_{a \in A} C_{FB}(a)$ *and* $C(a^*) < \infty \Rightarrow C'(a^*) > C(a^*)$.

PROOF: Let $(I_1', \ldots, I_n')$ be the cost minimizing way of implementing $a$ in the primed problem. Suppose that in the unprimed problem, the principal offers the agent the following *random* incentive scheme: for each $i$, if $q_i$ is the outcome, an $n$-sided die will be thrown where the probability of side $j$ coming up is $r_{ji}$, the $(j, i)$th element of $R$ ($j = 1, \ldots, n$). If side $j$ then comes up, you get $I_j'$. With this random incentive scheme, the probability of the agent getting $I_j'$ if he chooses a particular action is the same as in the primed problem. Therefore the agent's optimal action will be $a$. Furthermore, the principal's expected costs are the same as in the primed problem. This shows that the principal can implement $a$ at least as cheaply in the unprimed problem as in the primed problem by using a random incentive scheme. The final part of the proof is to note that the principal can reduce his expected cost further and continue to implement $a$ by offering the agent the perfectly certain utility level $v_i = \sum_{j=1}^n r_{ji} V(I_j')$ if the outcome is $q_i$ rather than the above lottery. That is, there is a deterministic incentive scheme which is better for the principal than the above random incentive scheme.

Q.E.D.

REMARK 8: The last part of the proof of Proposition 13 shows that it is never desirable under our assumptions for the principal to offer the agent an incentive scheme which makes his payment conditional on a particular outcome a lottery rather than a perfectly certain income.[17] This result may also be found in Holmstrom [7].

Note that if $\pi' = R\pi$ and $q'R = q$, the random variable $q'$ will have the same mean as $q$. In this case the following is true:

COROLLARY 1: *Make the hypotheses of Proposition* 13. *If, in addition,* $q'$ *is such that* $q'R = q$, *we have* $L' \geq L$.

PROOF: Obvious since $B'(a) = q'\pi'(a) = q'R\pi(a) = q\pi(a) = B(a)$.

---

[16]We use this notation to mean that every element of $R$ is strictly positive.

[17]This result depends strongly on our Assumption A1 that attitudes to income risk are independent of action. In the absence of this assumption, random incentive schemes may be desirable.

In the case $n = 2$, the transformation $\pi \to \pi' = R\pi$ is easy to interpret. Take any two actions $a_1, a_2 \in A$, and consider the likelihood ratio vector $(\pi_1(a_1)$ $/\pi_1(a_2), \pi_2(a_1)/\pi_2(a_2))$. Assume without loss of generality that $\pi_1(a_1)/\pi_1(a_2)$ $\le \pi_2(a_1)/\pi_2(a_2)$. Then it is easy to show that

$$(5.1) \qquad \left[ \frac{\pi_1'(a_1)}{\pi_1'(a_2)}, \frac{\pi_2'(a_1)}{\pi_2'(a_2)} \right] \subset \left[ \frac{\pi_1(a_1)}{\pi_1(a_2)}, \frac{\pi_2(a_1)}{\pi_2(a_2)} \right],$$

where $[x, y]$ is the interval between $x$ and $y$. In other words, the likelihood ratio vector becomes less variable in some sense when the stochastic transform $R$ is applied. In fact the converse to this is also true: if (5.1) holds, then there exists a stochastic matrix $R$ such that $\pi' = R\pi$ (see Blackwell and Girshick [3]). When $n > 2$, a simple characterization of this sort does not seem to exist, however.

One might ask whether a converse to Proposition 13 holds. That is, suppose $C'(a) \ge C(a)$ for all $a \in A$ and all concave utility functions $V$. Does it follow that $\pi'(a) = R\pi(a)$ for all $a \in A$, for some stochastic $R$? A converse along these lines can in fact be established when $n = 2$. Whether it holds for $n > 2$, we do not know.

Corollary 1 gives us a simple way of generating worse and worse incentive problems: repeatedly apply stochastic transforms to $\pi$. Suppose that we do this using always the same stochastic transform $R$, when $R \gg 0$ and is invertible. That is, we consider a sequence of incentive problems $1, 2, \ldots$, where in the $m$th problem $\pi_m(a) = R^{m-1}\pi(a)$ for all $a \in A$, and the gross profit vector $q_m$ satisfies $q_m R^{m-1} = q$ (this has a solution since $R$ is invertible). We know from Corollary 1 that $L_m$ will be increasing in $m$. The next proposition says that in the limit the loss from not being able to observe the agent reaches its maximal level.

DEFINITION: Let $L^* = \max_{a \in A}(B(a) - C_{FB}(a)) - \max\{B(a') - C_{FB}(a') \,|\, a'$ minimizes $C_{FB}(a)$ on $A\}$.

Since $C(a') = C_{FB}(a')$ if $a'$ minimizes $C_{FB}(a)$, $L^*$ is an upper limit on the loss to the principal from being unable to observe the agent. The next proposition shows that as the information $q$ reveals about $a$ gets smaller and smaller, the principal loses control over the agent, i.e., the agent chooses the least-cost action.

PROPOSITION 14: *Consider the sequence of incentive problems* $(A, U, \overline{U}, \pi_m, q_m)$, $m = 1, 2, \ldots$, *where* $\pi_m(a) = R^{m-1}\pi_1(a)$ *for all* $a \in A$, $q_m R^{m-1} = q_1$ *for some invertible stochastic matrix* $R \gg 0$. *Assume* A1, A2, *and* $\pi_{1i}(a) > 0$ *for all* $i = 1, \ldots, n$, *and* $a \in A$. *Then if* $V$ *is not a linear function,* $\lim_{m \to \infty} L_m = L^*$.

PROOF: It suffices to show that $\lim_{m \to \infty} C(a^*) = \infty$ for all $a^*$ with $C_{FB}(a^*)$ $> \min_{a \in A} C_{FB}(a)$. Suppose not for some such $a^*$. Let $(I_{m1}, \ldots, I_{mn})$ be the cost minimizing way of implementing $a^*$ in problem $m$. Then $\sum_i \pi_{mi}(a^*)I_{mi}$ and $\sum_i \pi_{mi}(a^*)V(I_{mi})$ are both bounded in $m$. It follows from Bertsekas [2] that the $(I_{mi})$ are bounded. Hence without loss of generality we may assume $I_{mi} \to I_i$ for

each $i$. It is easy to show that, since $R$ is a strictly positive stochastic matrix, $\lim_{m \to \infty} R^{m-1} = R^*$ where $R^*$ has the property that all of its columns are the same. Therefore $\lim_{m \to \infty} \pi_m(a) = R^* \pi_1(a) = \bar{\pi}$ is independent of $a$. But this means $\lim_{m \to \infty} \sum_i \pi_{mi}(a^*) V(I_{mi}) = \sum_i \bar{\pi}_i V(I_i) = \lim_{m \to \infty} \sum_i \pi_{mi}(a) V(I_{mi})$ for all $a \in A$. Hence the agent will prefer actions $a$ with $C_{FB}(a) < C_{FB}(a^*)$ to $a^*$. This contradicts the assumption that the incentive scheme implements $a^*$.        Q.E.D.

We turn now to a consideration of another factor which influences $L$: the agent's degree of risk aversion. Since no incentive problem arises when the agent is risk netural, but an incentive problem does arise when the agent is risk averse, one is led to ask whether $L$ increases as the agent becomes more risk averse. One difficulty in answering this question in general is the following. The way one makes the agent more risk averse is to replace his utility function $U(I, a)$ by $H(U(I, a))$ where $H$ is a real-valued, increasing, concave function. However, if $U$ satisfies Assumption A1, then $H(U)$ will generally not. To get around this difficulty, we will confine our attention to the case where $A$ is a subset of the real line, $V(I) = -e^{-kI}$, $G(a) = 0$, and $K(a) = e^{ka}$, i.e., the agent's utility function is $U(a, I) = -e^{-k(I-a)}$, where $k > 0$. Assume also that $\bar{U} = -e^{-k\alpha}$, i.e., the agent's outside opportunity is represented by the perfectly certain income $\alpha$. An increase in risk aversion can then be represented simply by an increase in $k$.

Note that if the agent's utility function is $-e^{-k(I-a)}$ and $\bar{U} = -e^{-k\alpha}$, then $C_{FB}(a) = a + \alpha$, which is independent of $k$. Hence first best profits are independent of $k$.

PROPOSITION 15: *Consider the incentive problem* $(A, U, \bar{U}, \pi, q)$ *where* $A$ *is a subset of the real line,* $U(a, I) = -e^{-k(I-a)}$, $\bar{U} = -e^{-k\alpha}$, *and* $k > 0$. *Assume* A3. *Write the loss from being unable to observe the agent as* $L(k)$. *Then* $\lim_{k \to 0} L(k) = 0$, $\lim_{k \to \infty} L(k) = L^*$.

PROOF: To show that $\lim_{k \to \infty} L(k) = L^*$, it suffices to show that $\lim_{k \to \infty} C(a^*, k) = \infty$ for all $a^*$ with $C_{FB}(a^*) > \min_{a \in A} C_{FB}(a)$. Suppose not for some such $a^*$, and let $C_{FB}(a) < C_{FB}(a^*)$. Then if $(I_1, \ldots, I_n)$ implements $a^*$, we must have

$$-\left(\sum_i \pi_i(a^*) e^{-kI_i}\right) e^{ka^*} \geq -\left(\sum_i \pi_i(a) e^{-kI_i}\right) e^{ka}$$

$(I_1, \ldots, I_n$ of course depend on $k$). Therefore,

$$(5.2) \qquad e^{k(a^* - a)} \leq \sum_i \pi_i(a) e^{-kI_i} \Big/ \sum_i \pi_i(a^*) e^{-kI_i}.$$

Now let $k \to \infty$. The LHS of (5.2) $\to \infty$. Therefore so must the RHS. We may assume w.l.o.g., however, that $I_1 = \min_i I_i$. Then

$$\frac{\sum_i \pi_i(a) e^{-kI_i}}{\sum_i \pi_i(a^*) e^{-kI_i}} = \frac{\sum_i \pi_i(a) e^{k(I_1 - I_i)}}{\sum_i \pi_i(a^*) e^{k(I_1 - I_i)}},$$

which is bounded since the denominator $\geq \pi_1(a^*)$. Contradiction.

We show now that $\lim_{k \to 0} L(k) = 0$. Let $I_i = q_i - F$. Then the agent maximizes

$$(5.3) \qquad E(-e^{-k(I-a)}) = -E\left(1 - k(I-a) + \frac{k^2}{2}(I-a)^2 + \cdots \right)$$

$$= -1 + k\left(\sum \pi_i(a)q_i - F - a\right) - \frac{k^2}{2}E(I-a)^2 + \cdots.$$

It follows that the agent maximizes

$$\left(\sum \pi_i(a)q_i - F - a\right) - \frac{k}{2}E(I-a)^2 + \cdots,$$

which means that in the limit $k \to 0$ the agent maximizes $B(a) - C_{FB}(a)$, i.e. chooses a first-best action. Furthermore, setting (5.3) equal to $-e^{-k\alpha} = -1 + k\alpha + \cdots$, we see that in the limit $k \to 0$,

$$\max_{a \in A}\left(\sum_i \pi_i(a)q_i - a\right) - F = \alpha,$$

so that the principal's expected profit equals $F = \max_{a \in A}(\sum_i \pi_i q_i - a) - \alpha = \max_{a \in A}(B(a) - C_{FB}(a)) = $ first-best profit.                     *Q.E.D.*

Proposition 15 tells us about the behavior of $L(k)$ for extreme values of $k$. It would be interesting to know whether $L(k)$ is increasing in $k$. We do not know the answer to that question except for the case $n = 2$, $A$ finite.

PROPOSITION 16: *Make the same hypotheses as in Proposition 14. Assume in addition that $n = 2$ and $A$ is finite. Then $L(k)$ is increasing in $k$.*

PROOF: See Appendix.

REMARK 9: Propositions 15 and 16 tell us how the principal's welfare varies with $k$. It is also interesting to ask how the shape of the optimal incentive scheme depends on $k$. Unfortunately, even in the case $n = 2$, very little can be said. In this case, the incentive scheme is characterized by the agent's share $s$. It is not difficult to construct examples showing that an increase in the agent's risk aversion may increase the optimal value of $s$, or may decrease it.

We conclude this section by considering how $L$ depends on the agent's incremental costs. Consider the case of additive separability, i.e., $K(a) \equiv$ constant. Suppose that we write the agent's utility function as $U_\lambda(a, I) = G_\lambda(a) + V(I)$, where $G_\lambda(a) = \alpha + \lambda F(a)$, $\lambda > 0$. (Without loss of generality, we take $K = 1$.) Then, when $\lambda$ is small, one feels that $L$ will be small since the agent does not require much of a reward to work hard. The fact that $\lim_{\lambda \to 0} L(\lambda) = 0$ has in fact been established by Shavell [20]. We prove a somewhat stronger result.

PROPOSITION 17: *Consider the incentive problem $(A, U_\lambda, \overline{U}, \pi, q)$, where $U_\lambda(a, I) = \alpha + \lambda F(a) + V(I)$ for all $a \in A$, $\lambda > 0$. Assume that A1–A3 hold for this*

*problem. Assume also that* (1) *A is an interval of the real line;* (2) *B(a) and F(a) are twice differentiable in the interior of A;* (3) *V is twice differentiable on $\mathscr{G}$ and $V' > 0$;* (4) *There is a unique maximizer $a^*$ of $B(a)$ lying in the interior of A and $B''(a^*) < 0$. Then $\lim_{\lambda \to 0}(L(\lambda)/\lambda) = 0$.*

PROOF: Consider the incentive problem with $\lambda = 1$. Then there are $a$'s arbitrarily close to $a^*$ for which $C(a)$ is finite. For let the principal set $v_i = rq_i - k$ where $k$ is chosen so that $v_i \in \mathscr{U}$ for all $i$. Then the agent will maximize $\sum \pi_i(a) U_\lambda(a, I_i)$, i.e. $\sum \pi_i(a) q_i + F(a)/r$. By letting $r \to \infty$, we can get the agent to choose an action arbitrarily close to $a^*$. For such an action, $C(a)$ will be finite.

Consider now an $a$ arbitrarily close to $a^*$. Let $(v_1, \ldots, v_n)$ be the cost minimizing way of implementing $a$ when $\lambda = 1$. Then it is clear from (2.2) that $(\lambda v_1 + \beta, \ldots, \lambda v_n + \beta)$ will implement $a$ for $\lambda \neq 1$, where

$$\lambda\left(\sum \pi_i(a) v_i + F(a)\right) + \alpha + \beta = \overline{U}.$$

It follows that

$$L(\lambda) \leq \sum \pi_i(\hat{a}) q_i - h\left(\overline{U} - \alpha - \lambda F(\hat{a})\right)$$

$$- \left(\sum \pi_i(a) q_i - \sum \pi_i(a) h(\lambda v_i + \beta)\right),$$

where $\hat{a}$ maximizes $\sum \pi_i(a) q_i - h(\overline{U} - \alpha - \lambda F(a))$, i.e. $\hat{a}$ is the first-best action in problem $\lambda$.

Therefore,

$$\frac{L(\lambda)}{\lambda} \leq \left[\frac{1}{\lambda}\left\{\sum \pi_i(a^*) q_i - h\left(\overline{U} - \alpha - \lambda F(a^*)\right)\right.\right.$$

$$\left.\left. - \left(\sum \pi_i(a) q_i - \sum \pi_i(a) h(\lambda v_i + \beta)\right)\right\}\right]$$

$$+ \left[\frac{1}{\lambda}\left\{\sum \pi_i(\hat{a}) q_i - h\left(\overline{U} - \alpha - \lambda F(\hat{a})\right)\right.\right.$$

$$\left.\left. - \sum \pi_i(a^*) q_i + h\left(\overline{U} - \alpha - \lambda F(a^*)\right)\right\}\right].$$

Now $\hat{a} \to a^*$ as $\lambda \to 0$. Furthermore, by differentiating the first-order conditions $(d/da)(\sum \pi_i(\hat{a}) q_i - h(\overline{U} - \alpha - \lambda F(\hat{a}))) = 0$, one can show that $d\hat{a}/d\lambda$ exists at $\lambda = 0$. It follows from the mean-value theorem and the fact that $B'(a^*) = 0$ that the second square bracket $\to 0$ as $\lambda \to 0$. To see that the first square bracket $\to 0$, note that, since $a$ is arbitrary, we can make $a$ converge to $a^*$ as fast as we like. Therefore we need only show that

(5.4)     $$\lim_{\lambda \to 0} \frac{1}{\lambda}\left(\sum \pi_i(a) h(\lambda v_i + \beta) - h\left(\overline{U} - \alpha - \lambda F(a^*)\right)\right) = 0.$$

But

$$\sum \pi_i(a) \left[ h(\lambda v_i + \beta) - h\left(\overline{U} - \alpha - \lambda F(a^*)\right) \right]$$

$$= \sum \pi_i(a) \left[ h\left(\lambda v_i + \overline{U} - \alpha - \lambda \sum \pi_j(a) v_j - \lambda F(a)\right) \right.$$

$$\left. - h\left(\overline{U} - \alpha - \lambda F(a^*)\right) \right]$$

$$= \sum \pi_i(a) \left[ h(\overline{U} - \alpha) + h'(\overline{U} - \alpha)\left(\lambda v_i - \lambda \sum \pi_j(a) v_j - \lambda F(a)\right) \right.$$

$$\left. + \cdots - h(\overline{U} - \alpha) + h'(\overline{U} - \alpha)(\lambda F(a^*)) + \cdots \right]$$

$$= h'(\overline{U} + \alpha)(-\lambda F(a) + \lambda F(a^*)) + \cdots$$

from which (5.4) follows. $Q.E.D.$

The proof of Proposition 17 is based on an envelope argument. It appears that a similar result can be established for the more general case where $U$ is not additively separable, but Assumption A1 holds. Since the proof is more complicated, however, we will not pursue this result here. The assumption that $a^*$ lies in the interior of $A$ may seem quite strong. Note, however, that if $a^*$ is a boundary point and $B'(a^*) \neq 0$, then the second-best optimal action equals $a^*$ for small enough $\lambda$. It is straightforward to apply the proof of Proposition 17 to show that $\lim_{\lambda \to 0}(L(\lambda)/\lambda) = 0$ in this case too.

Since the marginal product of labor of the agent—that is, the increase in expected profit resulting from an extra pound of expenditure by the agent—is proportional to $1/\lambda$, Proposition 17 can be interpreted as saying that the welfare loss $L$ is of a smaller order of magnitude than the reciprocal of the agent's marginal product of labor.

## 6. EXTENSIONS

We have assumed throughout the paper that the principal is risk-neutral and that the agent's attitudes to risk over income lotteries are independent of action —Assumption A1. We now briefly consider what happens if we relax these assumptions.

As we have noted in Section 2, Remark 3, our method of analysis generalizes without any difficulty to the case where the principal is risk-averse. Specific results change, however, The main difference is that now, even in the first-best situation, the principal will not bear all the risk. One implication of this is that even if there is no disutility of action for the agent, i.e. $a$ does not enter the agent's utility function, the first-best will not generally be reached. The reason is that there may be a conflict between the principal and agent over what income lottery should be selected (for a study of this conflict, see Ross [17] and Wilson [23]).

As a result of this, Proposition 3, part (5), is no longer true when the principal is risk-averse. Nor is Proposition 17 since $L(0) \neq 0$. Propositions 1 and 2 and Proposition 3, parts 1–4, continue to hold, however. So do Propositions 4 and 5 on the characterization of an optimal incentive scheme. Propositions 7, 8 generalize, as do Propositions 10, 11, and 12 (note that the function $C_{FB}$ is still well defined although it no longer refers to first-best cost). Proposition 3(6) does not hold and neither does Proposition 6 nor Proposition 9 (at least in its present form). Finally Corollary 1 of Proposition 13 and Propositions 14–16 do not generalize in an obvious way, since changing the risk aversion of the agent or the probability distribution of outcomes affects the first-best as well as the second-best.

The computational procedure presented in Section 4 for the two outcome case can be extended to the case where the principal is risk-averse. In the finite action case, it is still true that the agent will be indifferent between two actions at the optimum, except in the case where the first-best can be achieved. Thus it is necessary to check whether the first-best can be achieved. Otherwise the procedure is unaltered.

We turn now to the consequences of relaxing Assumption A1. These are much more serious since most of our analysis has depended crucially on being able to choose the control variables $V(I_1), \ldots, V(I_n)$ independently of $a$. Some results do generalize, however. In particular one can show that Propositions 1, 3, 10, and 12 generalize. It seems unlikely that the characterization of an optimal incentive scheme in Proposition 4 and Proposition 5, part 1, holds, but we do not have a counterexample. Surprisingly, perhaps, Proposition 5, part 2 does hold. Proposition 6 does not hold and it seems unlikely that Propositions 7–9 do.

In the two outcome case, one can still show that it is optimal for the agent's share $s$ to satisfy $0 \leq s < 1$. As a consequence Propositions 10 and 12 generalize. Proposition 11 does not generalize, however, and nor does our computational procedure for the two outcome case. Propositions 13 and 14 and Corollary 1 of Proposition 14 do not hold as they stand, although they do if one enlarges the set of feasible incentive schemes to include random schemes. (As we have noted in footnote 17, once Assumption A1 is dropped, random incentive schemes may be superior to deterministic schemes.) Finally, it seems likely that Proposition 17 could be generalized to the nonseparable case.

## 7. SUMMARY

The purpose of this paper has been to develop a method for analyzing the principal-agent problem in the case where the agent's attitudes to income risk are independent of action. Our method consists of breaking up the principal's problem into a computation of the costs and benefits accruing to the principal when the agent takes a particular action. We have used this method to establish a number of results about the structure of the optimal incentive scheme and about the determinants of the welfare loss resulting from the principal's inability to observe the agent's action. We have shown that it is never optimal for the

incentive scheme to be such that the principal's and agent's payoff are negatively related over the whole outcome range, although such a relationship may be optimal over part of the range. We have found sufficient conditions for the incentive scheme to be monotonic, progressive, and regressive. We have shown that a decrease in the quality of the principal's information in the sense of Blackwell increases welfare loss. When there are only two outcomes, welfare loss also increases when the agent becomes more risk averse. Finally, we have discussed how our techniques can be used to compute optimal incentive schemes in particular cases.

While we have talked throughout about "the" principal-agent problem, we have in fact been considering the simplest of a number of such problems. More complicated principal-agent problems arise when not only is the principal unable to monitor the agent, but also the agent possesses information about his environment, i.e. about $A$, $\pi$, or $U(a, I)$, which the principal does not. Such problems possess a number of features of the preference revelation problems studied in the recent incentive compatibility literature; see, for example, the *Review of Economic Studies* Symposium [16]. A start has been made in the analysis of such problems by Harris and Raviv [6], Holmstrom [7], and Mirrlees [12]. It will be interesting to see whether the techniques presented here will also be useful in the solution of these more complicated principal-agent problems.

*University of Chicago*
*and*
*London School of Economics*

### APPENDIX

PROOF OF PROPOSITION 16: It suffices to show that $C(a, k)$ is increasing locally in $k$ for each $a \in A$ whenever $C(a, k)$ is finite. Let $\tilde{k} = \lambda k$ $\lambda \geq 1$. Assume that $(I_1, I_2)$ is the cost minimizing way of implementing $a$, given $\tilde{k}$. Then, by the results of Section 4, e.g. equation (4.4),

(A1)
$$\pi_1 w_1 + \pi_2 w_2 = \frac{1}{e^{\tilde{k}(a+\alpha)}},$$

$$\pi_1' w_1 + \pi_2' w_2 = \frac{1}{e^{\tilde{k}(a'+\alpha)}},$$

where $w_1 = e^{-\tilde{k}I_1}$, $w_2 = e^{-\tilde{k}I_2}$, $\pi_1 = \pi_1(a)$, $\pi_2 = \pi_2(a)$, $\pi_1' = \pi_1(a')$, $\pi_2' = \pi_2(a')$, $a' \in A$, $a' < a$. Furthermore we can pick $a'$ so that $a'$ is independent of $k$ for $\lambda$ close to 1.

Equations (A1) determine $w_1$ and $w_2$ for each value of $\tilde{k}$. The cost of implementing $a$, $C(a, \tilde{k})$, is then given by

(A2) $\qquad C(a, \tilde{k}) = \pi_1 I_1 + \pi_2 I_2 = -\frac{1}{\tilde{k}}(\pi_1 \log w_1 + \pi_2 \log w_2).$

Differentiating (A2) with respect to $\lambda$ we get

(A3) $\qquad \dfrac{\partial C(a, \lambda k)}{\partial \lambda}\bigg|_{\lambda=1} = \dfrac{1}{k}\left(\pi_1 \log w_1 + \pi_2 \log w_2 - \dfrac{\pi_1}{w_1}\dfrac{dw_1}{d\lambda} - \dfrac{\pi_2}{w_2}\dfrac{dw_2}{d\lambda}\right).$

Set $x = e^{-k(a+\alpha)}$, $y = e^{-k(a'+\alpha)}$ in (A1). Then $e^{-\tilde{k}(a+\alpha)} = x^\lambda$, $e^{-\tilde{k}(a'+\alpha)} = y^\lambda$. Hence

$$
\begin{aligned}
\pi_1 \frac{dw_1}{d\lambda} + \pi_2 \frac{dw_2}{d\lambda} &= x \log x, \\
\pi_1' \frac{dw_1}{d\lambda} + \pi_2' \frac{dw_2}{d\lambda} &= y \log y,
\end{aligned}
$$
(A4)

where derivatives are evaluated at $\lambda = 1$. Solving (A1), (A4) yields

$$
w_1 = \frac{\pi_2' x - \pi_2 y}{\pi_1 \pi_2' - \pi_1' \pi_2} = \frac{\pi_2' x - \pi_2 y}{\pi_2' - \pi_2},
$$

$$
\frac{dw_1}{d\lambda} = \frac{\pi_2' x \log x - \pi_2 y \log y}{\pi_2' - \pi_2}.
$$

It follows that $\log w_1 \geq (1/w_1)(dw_1/d\lambda)$. For

$$
(A5) \quad w_1 \log w_1 - \frac{dw_1}{d\lambda} = \frac{\pi_2' x - \pi_2 y}{\pi_2' - \pi_2} \log \frac{\pi_2' x - \pi_2 y}{\pi_2' - \pi_2} - \left( \frac{\pi_2' x \log x - \pi_2 y \log y}{\pi_2' - \pi_2} \right)
$$

$$
= \frac{1}{\pi_2' - \pi_2} \left[ (\alpha x - \beta y) \log \frac{\alpha x - \beta y}{\alpha - \beta} - \alpha x \log x - \beta y \log y \right],
$$

where $\alpha = \pi_2'$, $\beta = \pi_2$. However, the RHS of (A5) $\geq 0$ by Lemma 3 below. The same argument shows that $\log w_2 \geq (1/w_2)(dw_2/d\lambda)$. It follows from (A3) that $(\partial C/\partial \lambda) \geq 0$, i.e., $C$ is increasing locally in $k$.

LEMMA 3: *Assume $\alpha, \beta, x, y > 0$. Then if $\alpha > \beta$ and $\alpha x > \beta y$, $\alpha x \log x - \beta y \log < (\alpha x - \beta y) \log((\alpha x - \beta y)/(\alpha - \beta))$. On the other hand, if $\alpha < \beta$ and $\alpha x < \beta y$, $\alpha x \log x - \beta y \log y > (\alpha x - \beta y) \log((\alpha x - \beta y)/(\alpha - \beta))$.*

PROOF: Since $z \log z$ is a convex function,

$$
\frac{\beta}{\alpha} (y \log y) + \left( \frac{\alpha - \beta}{\alpha} \right) \left( \frac{\alpha x - \beta y}{\alpha - \beta} \log \frac{\alpha x - \beta y}{\alpha - \beta} \right) \geq x \log x.
$$

This proves the first part. The second part follows similarly.                Q.E.D.

## REFERENCES

[1] ARROW, K. J.: "Insurance, Risk and Resource Allocation," *Essays in the Theory of Risk Bearing*. Chicago: Markham, 1971.
[2] BERTSEKAS, D.: "Necessary and Sufficient Conditions for Existence of an Optimal Portfolio," *Journal of Economic Theory*, 8(1974), 235–247.
[3] BLACKWELL, D., AND M. A. GIRSHICK: *Theory of Games and Statistical Decisions*. New York: John Wiley and Sons, Inc., 1954.
[4] BORCH, K.: *The Economics of Uncertainty*. Princeton: Princeton University Press, 1968.
[5] HARDY, G. H., J. E. LITTLEWOOD, AND G. POLYA: *Inequalities*. Cambridge: Cambridge University Press, 1952.
[6] HARRIS, M., AND A. RAVIV: "Optimal Incentive Contracts with Imperfect Information," *Journal of Economic Theory*, 20(1979), 231–259.
[7] HOLMSTROM, B.: "Moral Hazard and Observability," *Bell Journal of Economics*, 10(1979), 74–91.
[8] KEENEY, R.: "Risk Independence and Multiattributed Utility Functions," *Econometrica*, 41(1973), 27–34.
[9] MILGROM, P. R.: "Good News and Bad News: Representation Theorems and Applications," Discussion Paper No. 407, Northwestern University, Illinois, Mimeo, 1979.

[10] MIRRLEES, J. A.: "The Theory of Moral Hazard and Unobservable Behavior—Part I," Nuffield College, Oxford, Mimeo, 1975.

[11] ———: "The Optimal Structure of Incentives and Authority Within an Organization," *Bell Journal of Economics*, 7(1976), 105–131.

[12] ———: "The Implications of Moral Hazard for Optimal Insurance," Seminar given at Conference held in honour of Karl Borch, Bergen, Norway, Mimeo, 1979.

[13] PAULY, M.: "The Economics of Moral Hazard: Comment," *American Economic Review*, 58(1968), 531–536.

[14] POLLAK, R.: "The Risk Independence Axiom," *Econometrica*, 41(1973), 35–39.

[15] RADNER, R.: "Monitoring Cooperative Agreements in a Repeated Principal-Agent Relationship," Mimeo, Bell Laboratories, 1980.

[16] *Review of Economic Studies* Symposium on Incentive Compatibility, April, 1979.

[17] ROSS, S.: "The Economic Theory of Agency: The Principal's Problem," *American Economic Review*, 63(1973), 134–139.

[18] RUBINSTEIN, A. AND M. YAARI: Seminar given at Conference held in honour of Karl Borch, Bergen, Norway, 1979.

[19] SHAVELL, S.: "On Moral Hazard and Insurance," *Quarterly Journal of Economics*, 93(1979), 541–562.

[20] ———: "Risk Sharing and Incentives in the Principal and Agent Relationship," *Bell Journal of Economics*, 10(1979), 55–73.

[21] SPENCE, M., AND R. ZECKHAUSER: "Insurance, Information, and Individual Action," *American Economic Review*, 61(1971), 380–387.

[22] STIGLITZ, J. E.: "Incentives and Risk Sharing in Sharecropping," *Review of Economic Studies*, 61(1974), 219–256.

[23] WILSON, R.: "The Theory of Syndicates," *Econometrica*, 36(1968), 119–132.

[24] ZECKHAUSER, R.: "Medical Insurance: A Case Study of the Trade-Off Between Risk Spreading and Appropriate Incentives," *Journal of Economic Theory*, 2(1970), 10–26.

# Agency Theory: An Assessment and Review

**KATHLEEN M. EISENHARDT**
Stanford University

*Agency theory is an important, yet controversial, theory. This paper reviews agency theory, its contributions to organization theory, and the extant empirical work and develops testable propositions. The conclusions are that agency theory (a) offers unique insight into information systems, outcome uncertainty, incentives, and risk and (b) is an empirically valid perspective, particularly when coupled with complementary perspectives. The principal recommendation is to incorporate an agency perspective in studies of the many problems having a cooperative structure.*

One day Deng Xiaoping decided to take his grandson to visit Mao. "Call me granduncle," Mao offered warmly. "Oh, I certainly couldn't do that, Chairman Mao," the awe-struck child replied. "Why don't you give him an apple?" suggested Deng. No sooner had Mao done so than the boy happily chirped, "Oh thank you, Granduncle." "You see," said Deng, "what incentives can achieve." ("Capitalism," 1984, p. 62)

Agency theory has been used by scholars in accounting (e.g., Demski & Feltham, 1978), economics (e.g., Spence & Zeckhauser, 1971), finance (e.g., Fama, 1980), marketing (e.g., Basu, Lal, Srinivasan, & Staelin, 1985), political science (e.g., Mitnick, 1986), organizational behavior (e.g., Eisenhardt, 1985, 1988; Kosnik, 1987), and sociology (e.g., Eccles, 1985; White, 1985). Yet, it is still surrounded by controversy. Its proponents argue that a revolution is at hand and that "the foundation for a powerful theory of organizations is being put into place" (Jensen, 1983, p. 324). Its detractors call it trivial, dehumanizing, and even "dangerous" (Perrow, 1986, p. 235).

Which is it: grand theory or great sham? The purposes of this paper are to describe agency theory and to indicate ways in which organizational researchers can use its insights. The paper is organized around four questions that are germane to organizational research. The first asks the deceptively simple question, What is agency theory? Often, the technical style, mathematics, and tautological reasoning of the agency literature can obscure the theory. Moreover, the agency literature is split into two camps (Jensen, 1983), leading to differences in interpretation. For example, Barney and Ouchi (1986) argued that agency theory emphasizes how capital markets can affect the firm, whereas other authors made no reference to capital markets at all (Anderson, 1985; Demski & Feltham, 1978; Eccles, 1985; Eisenhardt, 1985).

The second question is, What does agency theory contribute to organizational theory? Proponents such as Ross (1973, p. 134) argued that "examples of agency are universal." Yet other scholars such as Perrow (1986) claimed that agency theory addresses no clear problems, and Hirsch and Friedman (1986) called it excessively narrow, focusing only on stock price. For economists, long accustomed to treating the or-

57

ganization as a "black box" in the theory of the firm, agency theory may be revolutionary. Yet, for organizational scholars the worth of agency theory is not so obvious.

The third question is, Is agency theory empirically valid? The power of the empirical research on agency theory to explain organizational phenomena is important to assess, particularly in light of the criticism that agency theory is "hardly subject to empirical test since it rarely tries to explain actual events" (Perrow, 1986, p. 224). Perrow (1986) also criticized the theory for being unrealistically one-sided because of its neglect of potential exploitation of workers.

The final question is, What topics and contexts are fruitful for organizational researchers who use agency theory? Identifying how useful agency theory can be to organizational scholars requires understanding the situations in which the agency perspective can provide theoretical leverage.

The principal contributions of the paper are to present testable propositions, identify contributions of the theory to organizational thinking, and evaluate the extant empirical literature. The overall conclusion is that agency theory is a useful addition to organizational theory. The agency theory ideas on risk, outcome uncertainty, incentives, and information systems are novel contributions to organizational thinking, and the empirical evidence is supportive of the theory, particularly when coupled with complementary theoretical perspectives.

## Origins of Agency Theory

During the 1960s and early 1970s, economists explored risk sharing among individuals or groups (e.g., Arrow, 1971; Wilson, 1968). This literature described the risk-sharing problem as one that arises when cooperating parties have different attitudes toward risk. Agency theory broadened this risk-sharing literature to include the so-called agency problem that occurs when cooperating parties have different goals and di-

vision of labor (Jensen & Meckling, 1976; Ross, 1973). Specifically, agency theory is directed at the ubiquitous agency relationship, in which one party (the principal) delegates work to another (the agent), who performs that work. Agency theory attempts to describe this relationship using the metaphor of a contract (Jensen & Meckling, 1976).

Agency theory is concerned with resolving two problems that can occur in agency relationships. The first is the *agency problem* that arises when (a) the desires or goals of the principal and agent conflict and (b) it is difficult or expensive for the principal to verify what the agent is actually doing. The problem here is that the principal cannot verify that the agent has behaved appropriately. The second is the *problem of risk sharing* that arises when the principal and agent have different attitudes toward risk. The problem here is that the principal and the agent may prefer different actions because of the different risk preferences.

Because the unit of analysis is the contract governing the relationship between the principal and the agent, the focus of the theory is on determining the most efficient contract governing the principal-agent relationship given assumptions about people (e.g., self-interest, bounded rationality, risk aversion), organizations (e.g., goal conflict among members), and information (e.g., information is a commodity which can be purchased). Specifically, the question becomes, Is a behavior-oriented contract (e.g., salaries, hierarchical governance) more efficient than an outcome-oriented contract (e.g., commissions, stock options, transfer of property rights, market governance)? An overview of agency theory is given in Table 1.

The agency structure is applicable in a variety of settings, ranging from macrolevel issues such as regulatory policy to microlevel dyad phenomena such as blame, impression management, lying, and other expressions of self-interest. Most frequently, agency theory has been applied to organizational phenomena

58

**Table 1**
*Agency Theory Overview*

| | |
|---|---|
| Key idea | Principal-agent relationships should reflect efficient organization of information and risk-bearing costs |
| Unit of analysis | Contract between principal and agent |
| Human assumptions | Self-interest<br>Bounded rationality<br>Risk aversion |
| Organizational assumptions | Partial goal conflict among participants<br>Efficiency as the effectiveness criterion<br>Information asymmetry between principal and agent |
| Information assumption | Information as a purchasable commodity |
| Contracting problems | Agency (moral hazard and adverse selection)<br>Risk sharing |
| Problem domain | Relationships in which the principal and agent have partly differing goals and risk preferences (e.g., compensation, regulation, leadership, impression management, whistle-blowing, vertical integration, transfer pricing) |

such as compensation (e.g., Conlon & Parks, 1988; Eisenhardt, 1985), acquisition and diversification strategies (e.g., Amihud & Lev, 1981), board relationships (e.g., Fama & Jensen, 1983; Kosnik, 1987), ownership and financing structures (e.g., Argawal & Mandelker, 1987; Jensen & Meckling, 1976), vertical integration (Anderson, 1985; Eccles, 1985), and innovation (Bolton, 1988; Zenger, 1988). Overall, the domain of agency theory is relationships that mirror the basic agency structure of a principal and an agent who are engaged in cooperative behavior, but have differing goals and differing attitudes toward risk.

## Agency Theory

From its roots in information economics, agency theory has developed along two lines:

positivist and principal-agent (Jensen, 1983). The two streams share a common unit of analysis: the contract between the principal and the agent. They also share common assumptions about people, organizations, and information. However, they differ in their mathematical rigor, dependent variable, and style.

### Positivist Agency Theory

Positivist researchers have focused on identifying situations in which the principal and agent are likely to have conflicting goals and then describing the governance mechanisms that limit the agent's self-serving behavior. Positivist research is less mathematical than principal-agent research. Also, positivist researchers have focused almost exclusively on the special case of the principal-agent relationship between owners and managers of large, public corporations (Berle & Means, 1932).

Three articles have been particularly influential. Jensen and Meckling (1976) explored the ownership structure of the corporation, including how equity ownership by managers aligns managers' interests with those of owners. Fama (1980) discussed the role of efficient capital and labor markets as information mechanisms that are used to control the self-serving behavior of top executives. Fama and Jensen (1983) described the role of the board of directors as an information system that the stockholders within large corporations could use to monitor the opportunism of top executives. Jensen and his colleagues (Jensen, 1984; Jensen & Roeback, 1983) extended these ideas to controversial practices, such as golden parachutes and corporate raiding.

From a theoretical perspective, the positivist stream has been most concerned with describing the governance mechanisms that solve the agency problem. Jensen (1983, p. 326) described this interest as "why certain contractual relations arise." Two propositions capture the governance mechanisms which are identified in the positivist stream. One proposition is that out-

59

come-based contracts are effective in curbing agent opportunism. The argument is that such contracts coalign the preferences of agents with those of the principal because the rewards for both depend on the same actions, and, therefore, the conflicts of self-interest between principal and agent are reduced. For example, Jensen and Meckling (1976) described how increasing the firm ownership of the managers decreases managerial opportunism. In formal terms,

*Proposition 1: When the contract between the principal and agent is outcome based, the agent is more likely to behave in the interests of the principal.*

The second proposition is that information systems also curb agent opportunism. The argument here is that, since information systems inform the principal about what the agent is actually doing, they are likely to curb agent opportunism because the agent will realize that he or she cannot deceive the principal. For example, Fama (1980) described the information effects of efficient capital and labor markets on managerial opportunism, and Fama and Jensen (1983) described the information role that boards of directors play in controlling managerial behavior. In formal terms,

*Proposition 2: When the principal has information to verify agent behavior, the agent is more likely to behave in the interests of the principal.*

At its best, positivist agency theory can be regarded as enriching economics by offering a more complex view of organizations (Jensen, 1983). However, it has been criticized by organizational theorists as minimalist (Hirsch, Michaels, & Friedman, 1987; Perrow, 1986) and by microeconomists as tautological and lacking rigor (Jensen, 1983). Nonetheless, positivist agency theory has ignited considerable research (Barney & Ouchi, 1986) and popular interest ("Meet Mike," 1988).

## Principal-Agent Research

Principal-agent researchers are concerned with a general theory of the principal-agent relationship, a theory that can be applied to employer-employee, lawyer-client, buyer-supplier, and other agency relationships (Harris & Raviv, 1978). Characteristic of formal theory, the principal-agent paradigm involves careful specification of assumptions, which are followed by logical deduction and mathematical proof.

In comparison with the positivist stream, principal-agent theory is abstract and mathematical and, therefore, less accessible to organizational scholars. Indeed, the most vocal critics of the theory (Perrow, 1986; Hirsch et al., 1987) have focused their attacks primarily on the more widely known positivist stream. Also, the principal-agent stream has a broader focus and greater interest in general, theoretical implications. In contrast, the positivist writers have focused almost exclusively on the special case of the owner/CEO relationship in the large corporation. Finally, principal-agent research includes many more testable implications.

For organizational scholars, these differences provide background for understanding criticism of the theory. However, they are not crucial. Rather, the important point is that the two streams are complementary: Positivist theory identifies various contract alternatives, and principal-agent theory indicates which contract is the most efficient under varying levels of outcome uncertainty, risk aversion, information, and other variables described below.

The focus of the principal-agent literature is on determining the optimal contract, behavior versus outcome, between the principal and the agent. The simple model assumes goal conflict between principal and agent, an easily measured outcome, and an agent who is more risk averse than the principal. (Note: The argument behind a more risk averse agent is that agents, who are unable to diversify their employment, should be risk averse and principals, who are

60

185

capable of diversifying their investments, should be risk neutral.) The approach of the simple model can be described in terms of cases (e.g., Demski & Feltham, 1978). The first case, a simple case of complete information, is when the principal knows what the agent has done. Given that the principal is buying the agent's behavior, then a contract that is based on behavior is most efficient. An outcome-based contract would needlessly transfer risk to the agent, who is assumed to be more risk averse than the principal.

The second case is when the principal does not know exactly what the agent has done. Given the self-interest of the agent, the agent may or may not have behaved as agreed. The agency problem arises because (a) the principal and the agent have different goals and (b) the principal cannot determine if the agent has behaved appropriately. In the formal literature, two aspects of the agency problem are cited. *Moral hazard* refers to lack of effort on the part of the agent. The argument here is that the agent may simply not put forth the agreed-upon effort. That is, the agent is shirking. For example, moral hazard occurs when a research scientist works on a personal research project on company time, but the research is so complex that corporate management cannot detect what the scientist is actually doing. *Adverse selection* refers to the misrepresentation of ability by the agent. The argument here is that the agent may claim to have certain skills or abilities when he or she is hired. Adverse selection arises because the principal cannot completely verify these skills or abilities either at the time of hiring or while the agent is working. For example, adverse selection occurs when a research scientist claims to have experience in a scientific specialty and the employer cannot judge whether this is the case.

In the case of unobservable behavior (due to moral hazard or adverse selection), the principal has two options. One is to discover the agent's behavior by investing in information systems such as budgeting systems, reporting procedures, boards of directors, and additional layers of management. Such investments reveal the agent's behavior to the principal, and the situation reverts to the complete information case. In formal terms,

> *Proposition 3: Information systems are positively related to behavior-based contracts and negatively related to outcome-based contracts.*

The other option is to contract on the outcomes of the agent's behavior. Such an outcome-based contract motivates behavior by coalignment of the agent's preferences with those of the principal, but at the price of transferring risk to the agent. The issue of risk arises because outcomes are only partly a function of behaviors. Government policies, economic climate, competitor actions, technological change, and so on, may cause uncontrollable variations in outcomes. The resulting outcome uncertainty introduces not only the inability to preplan, but also risk that must be borne by someone. When outcome uncertainty is low, the costs of shifting risk to the agent are low and outcome-based contracts are attractive. However, as uncertainty increases, it becomes increasingly expensive to shift risk despite the motivational benefits of outcome-based contracts. In formal terms,

> *Proposition 4: Outcome uncertainty is positively related to behavior-based contracts and negatively related to outcome-based contracts.*

This simple agency model has been described in varying ways by many authors (e.g., Demski & Feltham, 1978; Harris & Raviv, 1979; Holmstrom, 1979; Shavell, 1979). However, the heart of principal-agent theory is the trade-off between (a) the cost of measuring behavior and (b) the cost of measuring outcomes and transferring risk to the agent.

A number of extensions to this simple model are possible. One is to relax the assumption of a risk-averse agent (e.g., Harris & Raviv, 1979). Research (MacCrimmon & Wehrung, 1986) indicates that individuals vary widely in their risk

61

186

attitudes. As the agent becomes increasingly less risk averse (e.g., a wealthy agent), it becomes more attractive to pass risk to the agent using an outcome-based contract. Conversely, as the agent becomes more risk averse, it is increasingly expensive to pass risk to the agent. In formal terms,

> Proposition 5: The risk aversion of the agent is positively related to behavior-based contracts and negatively related to outcome-based contracts.

Similarly, as the principal becomes more risk averse, it is increasingly attractive to pass risk to the agent. In formal terms,

> Proposition 6: The risk aversion of the principal is negatively related to behavior-based contracts and positively related to outcome-based contracts.

Another extension is to relax the assumption of goal conflict between the principal and agent (e.g., Demski, 1980). This might occur either in a highly socialized or clan-oriented firm (Ouchi, 1979) or in situations in which self-interest gives way to selfless behavior (Perrow, 1986). If there is no goal conflict, the agent will behave as the principal would like, regardless of whether his or her behavior is monitored. As goal conflict decreases, there is a decreasing motivational imperative for outcome-based contracting, and the issue reduces to risk-sharing considerations. Under the assumption of a risk-averse agent, behavior-based contracts become more attractive. In formal terms,

> Proposition 7: The goal conflict between principal and agent is negatively related to behavior-based contracts and positively related to outcome-based contracts.

Another set of extensions relates to the task performed by the agent. For example, the progammability of the task is likely to influence the ease of measuring behavior (Eisenhardt, 1985, 1988). *Programmability* is defined as the degree to which appropriate behavior by the agent can be specified in advance. For example, the job of a retail sales cashier is much more programmed than that of a high-technology entrepreneur. The argument is that the behavior of agents engaged in more programmed jobs is easier to observe and evaluate. Therefore, the more programmed the task, the more attractive are behavior-based contracts because information about the agent's behavior is more readily determined. Very programmed tasks readily reveal agent behavior, and the situation reverts to the complete information case. Thus, retail sales clerks are more likely to be paid via behavior-based contracting (e.g., hourly wages), whereas entrepreneurs are more likely to be compensated with outcome-based contracts (e.g., stock ownership). In formal terms,

> Proposition 8: Task programmability is positively related to behavior-based contracts and negatively related to outcome-based contracts.

Another task characteristic is the measurability of the outcome (Anderson, 1985; Eisenhardt, 1985). The simple model assumes that outcomes are easily measured. However, some tasks require a long time to complete, involve joint or team effort, or produce soft outcomes. In these circumstances, outcomes are either difficult to measure or difficult to measure within a practical amount of time. When outcomes are measured with difficulty, outcome-based contracts are less attractive. In contrast, when outcomes are readily measured, outcome-based contracts are more attractive. In formal terms,

> Proposition 9: Outcome measurability is negatively related to behavior-based contracts and positively related to outcome-based contracts.

Finally, it seems reasonable that when principals and agents engage in a long-term relationship, it is likely that the principal will learn about the agent (e.g., Lambert, 1983) and so will be able to assess behavior more readily. Conversely, in short-term agency relationships, the information asymmetry between principal and agent is likely to be greater, thus making out-

62

come-based contracts more attractive. In formal terms,

*Proposition 10: The length of the agency relationship is positively related to behavior-based contracts and negatively related to outcome-based contracts.*

## Agency Theory and the Organizational Literature

Despite Perrow's (1986) assertion that agency theory is very different from organization theory, agency theory has several links to mainstream organization perspectives (see Table 2). At its roots, agency theory is consistent with the classic works of Barnard (1938) on the nature of co-operative behavior and March and Simon (1958) on the inducements and contributions of the employment relationship. As in this earlier work, the heart of agency theory is the goal conflict inherent when individuals with differing preferences engage in cooperative effort, and the essential metaphor is that of the contract.

Agency theory is also similar to political models of organizations. Both agency and political perspectives assume the pursuit of self-interest at the individual level and goal conflict at the organizational level (e.g., March, 1962; Pfeffer, 1981). Also, in both perspectives, information

asymmetry is linked to the power of lower order participants (e.g., Pettigrew, 1973). The difference is that in political models goal conflicts are resolved through bargaining, negotiation, and coalitions—the power mechanism of political science. In agency theory they are resolved through the coalignment of incentives—the price mechanism of economics.

Agency theory also is similar to the information processing approaches to contingency theory (Chandler, 1962; Galbraith, 1973; Lawrence & Lorsch, 1967). Both perspectives are information theories. They assume that individuals are boundedly rational and that information is distributed asymmetrically throughout the organization. They also are efficiency theories; that is, they use efficient processing of information as a criterion for choosing among various organizing forms (Galbraith, 1973). The difference between the two is their focus: In contingency theory researchers are concerned with the optimal structuring of reporting relationships and decision-making responsibilities (e.g., Galbraith, 1973; Lawrence & Lorsch, 1967), whereas in agency theory they are concerned with the optimal structuring of control relationships resulting from these reporting and decision-making patterns. For example, using contingency theory, we would be concerned with whether a firm is organized in a divisional or matrix structure.

**Table 2**
***Comparison of Agency Theory Assumptions and Organizational Perspectives***

| | Perspective | | | | |
| --- | --- | --- | --- | --- | --- |
| **Assumption** | **Political** | **Contingency** | **Organization Control** | **Transaction Cost** | **Agency** |
| Self-interest | X | | | X | X |
| Goal conflict | X | | | X | X |
| Bounded rationality | | X | X | X | X |
| Information asymmetry | | X | | X | X |
| Preeminence of efficiency | | X | X | X | X |
| Risk aversion | | | | | X |
| Information as a commodity | | | | | X |

63

Using agency theory, we would be concerned with whether managers within the chosen structure are compensated by performance incentives.

The most obvious tie is with the organizational control literature (e.g., Dornbusch & Scott, 1974). For example, Thompson's (1967) and later Ouchi's (1979) linking of known means/ends relationships and crystallized goals to behavior versus outcome control is very similar to agency theory's linking task programmability and measurability of outcomes to contract form (Eisenhardt, 1985). That is, known means/ends relationships (task programmability) lead to behavior control, and crystallized goals (measurable outcomes) lead to outcome control. Similarly, Ouchi's (1979) extension of Thompson's (1967) framework to include clan control is similar to assuming low goal conflict (Proposition 7) in agency theory. Clan control implies goal congruence between people and, therefore, the reduced need to monitor behavior or outcomes. Motivation issues disappear. The major differences between agency theory and the organizational control literature are the risk implications of principal and agent risk aversion and outcome uncertainty (Propositions 4, 5, 6).

Not surprisingly, agency theory has similarities with the transaction cost perspective (Williamson, 1975). As noted by Barney and Ouchi (1986), the theories share assumptions of self-interest and bounded rationality. They also have similar dependent variables; that is, hierarchies roughly correspond to behavior-based contracts, and markets correspond to outcome-based contracts. However, the two theories arise from different traditions in economics (Spence, 1975): In transaction cost theorizing we are concerned with organizational boundaries, whereas in agency theorizing the contract between cooperating parties, regardless of boundary, is highlighted. However, the most important difference is that each theory includes unique independent variables. In transaction cost theory these are asset specificity and small numbers bargaining. In agency theory there are the risk attitudes of the principal and agent, outcome uncertainty, and information systems. Thus, the two theories share a parentage in economics, but each has its own focus and several unique independent variables.

## Contributions of Agency Theory

Agency theory reestablishes the importance of incentives and self-interest in organizational thinking (Perrow, 1986). Agency theory reminds us that much of organizational life, whether we like it or not, is based on self-interest. Agency theory also emphasizes the importance of a common problem structure across research topics. As Barney and Ouchi (1986) described it, organization research has become increasingly topic, rather than theory, centered. Agency theory reminds us that common problem structures do exist across research domains. Therefore, results from one research area (e.g., vertical integration) may be germane to others with a common problem structure (e.g., compensation).

Agency theory also makes two specific contributions to organizational thinking. The first is the treatment of information. In agency theory, information is regarded as a commodity: It has a cost, and it can be purchased. This gives an important role to formal information systems, such as budgeting, MBO, and boards of directors, and informal ones, such as managerial supervision, which is unique in organizational research. The implication is that organizations can invest in information systems in order to control agent opportunism.

An illustration of this is executive compensation. A number of authors in this literature have expressed surprise at the lack of performance-based executive compensation (e.g., Pearce, Stevenson, & Perry, 1985; Ungson & Steers, 1984). However, from an agency perspective, it is not surprising since such compensation should be contingent upon a variety of factors including information systems. Specifically,

64

richer information systems control managerial opportunism and, therefore, lead to less performance-contingent pay.

One particularly relevant information system for monitoring executive behaviors is the board of directors. From an agency perspective, boards can be used as monitoring devices for shareholder interests (Fama & Jensen, 1983). When boards provide richer information, compensation is less likely to be based on firm performance. Rather, because the behaviors of top executives are better known, compensation based on knowledge of executive behaviors is more likely. Executives would then be rewarded for taking well-conceived actions (e.g., high risk/high potential R&D) whose outcomes may be unsuccessful. Also, when boards provide richer information, top executives are more likely to engage in behaviors that are consistent with stockholders' interests. For example, from an agency viewpoint, behaviors such as using greenmail and golden parachutes, which tend to benefit the manager more than the stockholders, are less likely when boards are better monitors of stockholders' interests. Operationally, the richness of board information can be measured in terms of characteristics such as frequency of board meetings, number of board subcommittees, number of board members with long tenure, number of board members with managerial and industry experience, and number of board members representing specific ownership groups.

A second contribution of agency theory is its risk implications. Organizations are assumed to have uncertain futures. The future may bring prosperity, bankruptcy, or some intermediate outcome, and that future is only partly controlled by organization members. Environmental effects such as government regulation, emergence of new competitors, and technical innovation can affect outcomes. Agency theory extends organizational thinking by pushing the ramifications of outcome uncertainty to their implications for creating risk. Uncertainty is viewed in terms of risk/reward trade-offs, not just in terms of inability to preplan. The implication is that outcome uncertainty coupled with differences in willingness to accept risk should influence contracts between principal and agent.

Vertical integration provides an illustration. For example, Walker and Weber (1984) found that technological and demand uncertainty did not affect the "make or buy" decision for components in a large automobile manufacturer (principal in this case). The authors were unable to explain their results using a transaction cost framework. However, their results are consistent with agency thinking if the managers of the automobile firm are risk neutral (a reasonable assumption given the size of the automobile firm relative to the importance of any single component). According to agency theory, we would predict that such a risk-neutral principal is relatively uninfluenced by outcome uncertainty, which was Walker and Weber's result.

Conversely, according to agency theory, the reverse prediction is true for a new venture. In this case, the firm is small and new, and it has limited resources available to it for weathering uncertainty: The likelihood of failure looms large. In this case, the managers of the venture may be risk-averse principals. If so, according to agency theory we would predict that such managers will be very sensitive to outcome uncertainty. In particular, the managers would be more likely to choose the "buy" option, thereby transferring risk to the supplying firm. Overall, agency theory predicts that risk-neutral managers are likely to choose the "make" option (behavior-based contract), whereas risk-averse executives are likely to choose "buy" (outcome-based contract).

## Empirical Results

Researchers in several disciplines have undertaken empirical studies of agency theory. These studies, mirroring the two streams of theoretical agency research, are in Table 3.

65

190

**Table 3**
**Summary of Agency Theory Studies**

| Author(s) | Research Stream | Sample | Agency Variables | Companion Theory | Dependent Variables | Results |
|---|---|---|---|---|---|---|
| Amihud & Lev (1981) | Positivist | 309 Fortune 500 firms | Manager vs. owner controlled | None | Conglomerate mergers & diversification | Support |
| Walking & Long (1984) | Positivist | 105 U.S. firms | Management's equity & options | Shareholder welfare & other controls | Managerial resistance to takeover bid | Support |
| Anderson (1985) | Principal-Agent | 159 sales districts in 13 electronics firms | Importance of nonselling activities, length of selling cycle, & difficulty evaluating sales performance | Transaction cost | Representative vs. corporate sales force | Mixed |
| Eisenhardt (1985) | Principal-Agent | 54 retail stores | Information systems, cost of outcome measurement, & outcome uncertainty | Organizational control | Salary vs. commission | Support |
| Eccles (1985) | Principal-Agent | 150 interviews in 13 chemical, electronics, heavy machinery, & machine component firms | Decentralization | Equity | Type of transfer price | Inductive model |
| Wolfson (1985) | Positivist | 39 oil & gas limited partnerships | General partner's track record | Tax effects | Share price | Support |
| Argawal & Mandelker (1987) | Positivist | 209 major corporations | Executive stock holdings | None | Acquisitions, divestitures, & debt/equity ratio | Support |

66

**Table 3 (continued)**
***Summary of Agency Theory Studies***

| Author(s) | Research Stream | Sample | Agency Variables | Companion Theory | Dependent Variables | Results |
|---|---|---|---|---|---|---|
| Kosnik (1987) | Positivist | 110 major corporations targeted for greenmail | Proportion of outside directors, equity held by outside directors, & outside directors with executive experience | Hegemony | Payment of greenmail (Yes/No) | Mixed |
| Eisenhardt (1988) | Principal-Agent | 54 retail stores | Job program-mability, span of control, & outcome uncertainty | Institutional | Salary vs. commission | Support |
| Conlon & Parks (1988) | Principal-Agent | 40 dyads | Monitoring | Institutional | Performance-contingent compensation | Support |
| Barney (1988) | Positivist | 32 Japanese electronics firms | Employee stock ownership | Size & growth controls | Cost of equity | Support |
| Singh & Harianto (in press) | Positivist | 84 *Fortune 500* firms | Managerial stock ownership & takeover threat | Managerialist | Golden parachute contracts | Support |

*Note.* This set of studies was developed through contacting other agency researchers, scanning journals, and following up referenced articles. Although the list is not exhaustive, it includes many of the relevant studies.

67

## Results of the Positivist Stream

In the positivist stream, the common approach is to identify a policy or behavior in which stockholder and management interests diverge and then to demonstrate that information systems or outcome-based incentives solve the agency problem. That is, these mechanisms coalign managerial behaviors with owner preferences. Consistent with the positivist tradition, most of these studies concern the separation of ownership from management in large corporations, and they use secondary source data that are available for large firms.

One of the earliest studies of this type was conducted by Amihud and Lev (1981). These researchers explored why firms engage in conglomerate mergers. In general, conglomerate mergers are not in the interests of the stockholders because, typically, stockholders can diversify directly through their stock portfolio. In contrast, conglomerate mergers may be attractive to managers who have fewer avenues available to diversify their own risk. Hence, conglomerate mergers are an arena in which owner and manager interests diverge. Specifically, these authors linked merger and diversification behaviors to whether the firm was owner controlled (i.e., had a major stockholder) or manager controlled (i.e., had no major stockholder). Consistent with agency theory arguments (Jensen & Meckling, 1976), manager-controlled firms engaged in significantly more conglomerate (but not more related) acquisitions and were more diversified.

Along the same lines, Walking and Long (1984) studied managers' resistance to takeover bids. Their sample included 105 large U.S. corporations that were targets of takeover attempts between 1972 and 1977. In general, resistance to takeover bids is not in the stockholders' interests, but it may be in the interests of managers because they can lose their jobs during a takeover. Consistent with agency theory (Jensen & Meckling, 1976), the authors found that managers who have substantial equity positions within

their firms (outcome-based contracts) were less likely to resist takeover bids.

The effects of market discipline on agency relationships were examined in Wolfson's (1985) study of the relationship between the limited (principals) and general (agent) partners in oil and gas tax shelter programs. In this study, both tax and agency effects were combined in order to assess why the limited partnership governance form survived in this setting despite extensive information advantages and divergent incentives for the limited partner. Consistent with agency arguments (Fama, 1980), Wolfson found that long-run reputation effects of the market coaligned the short-run behaviors of the general partner with the limited partners' welfare.

Kosnik (1987) examined another information mechanism for managerial opportunism, the board of directors. Kosnik studied 110 large U.S. corporations that were greenmail targets between 1979 and 1983. Using both hegemony and agency theories, she related board characteristics to whether greenmail was actually paid (paying greenmail is considered not in the stockholders' interests). As predicted by agency theory (Fama & Jensen, 1983), boards of companies that resisted greenmail had a higher proportion of outside directors and a higher proportion of outside director executives.

In a similar vein, Argawal and Mandelker (1987) examined whether executive holdings of firm securities reduced agency problems between stockholders and management. Specifically, they studied the relationship between stock and stock option holdings of executives and whether acquisition and financing decisions were made consistent with the interests of stockholders. In general, managers prefer lower risk acquisitions and lower debt financing (see Argawal & Mandelker, 1987, for a review). Their sample included 209 firms that participated in acquisitions and divestitures between 1974 and 1982. Consistent with agency ideas (e.g., Jensen & Meckling, 1976), executive security holdings (outcome-based contract) were related to acqui-

68

sition and financing decisions that were more consistent with stockholder interest. That is, executive stock holdings appeared to coalign managerial preferences with those of stockholders.

Singh and Harianto (in press) studied golden parachutes in a matched sample of 84 *Fortune* 500 firms. Their study included variables from both agency and managerialist perspectives. Consistent with agency theory (Jensen & Meckling, 1976; Fama & Jensen, 1983), the authors found that golden parachutes are used to coalign executive interests with those of stockholders in takeover situations, and they are seen as an alternative outcome-based contract to executive stock ownership. Specifically, the authors found that golden parachutes were positively associated with a higher probability of a takeover attempt and negatively associated with executive stock holdings.

Finally, Barney (1988) explored whether employee stock ownership reduces a firm's cost of equity capital. Consistent with agency theory (Jensen & Meckling, 1976), Barney argued that employee stock ownership (outcome-based contract) would coalign the interests of employees with stockholders. Using efficient capital market assumptions, he further argued that this coalignment would be reflected in the market through a lower cost of equity. Although Barney did not directly test the agency argument, the results are consistent with an agency view.

In summary, there is support for the existence of agency problems between shareholders and top executives across situations in which their interests diverge—that is, takeover attempts, debt versus equity financing, acquisitions, and divestitures, and for the mitigation of agency problems (a) through outcome-based contracts such as golden parachutes (Singh & Harianto, in press) and executive stock holdings (Argawal & Mandelker, 1987; Walking & Long, 1984) and (b) through information systems such as boards (Kosnik, 1987) and efficient markets (Barney, 1988; Wolfson, 1985). Overall, these studies sup-

port the positivist propositions described earlier. Similarly, laboratory studies by DeJong and colleagues (1985), which are not reviewed here, are also supportive.

## Results of the Principal-Agent Stream

The principal-agent stream is more directly focused on the contract between the principal and the agent. Whereas the positivist stream lays the foundation (that is, that agency problems exist and that various contract alternatives are available), the principal-agent stream indicates the most efficient contract alternative in a given situation. The common approach in these studies is to use a subset of agency variables such as task programmability, information systems, and outcome uncertainty to predict whether the contract is behavior- or outcome-based. The underlying assumption is that principals and agents will choose the most efficient contract, although efficiency is not directly tested.

In one study, Anderson (1985) probed vertical integration using a transaction cost perspective with agency variables. Specifically, she examined the choice between a manufacturer's representative (outcome-based) and a corporate sales force (behavior-based) among a sample of electronics firms. The most powerful explanatory variable was from agency theory: the difficulty of measuring outcomes (measured by amount of nonselling tasks and joint team sales). Consistent with agency predictions, this variable was positively related to using a corporate sales force (behavior-based contract).

In other studies, Eisenhardt (1985, 1988) examined the choice between commission (outcome-based) and salary (behavior-based) compensation of salespeople in retailing. The original study (1985) included only agency variables, while a later study (1988) added additional agency variables and institutional theory predictions. The results supported agency theory predictions that task programmability, information systems (measured by the span of control), and outcome uncertainty variables (measured

69

194

by number of competitors and failure rates) sig-nificantly predict the salary versus commission choice. Institutional variables were significant as well.

Conlon and Parks (1988) replicated and ex-tended Eisenhardt's work in a laboratory set-ting. They used a multiperiod design to test both agency and institutional predictions. Consistent with agency theory (Harris & Raviv, 1978), they found that information systems (manipulated by whether or not the principal could monitor the agent's behavior) were negatively related to performance-contingent (outcome-based) pay. They also found support for the institutional pre-dictions.

Finally, Eccles (1985) used agency theory to develop a framework for understanding transfer pricing. Using interviews with 150 executives in 13 large corporations, he developed a frame-work based on notions of agency and fairness to prescribe the conditions under which various sourcing and transfer pricing alternatives are both efficient and equitable. Prominent in his framework is the link between decentralization (arguably a measure of task programmability) and the choice between cost (behavior-based contract) and market (outcome-based contract) transfer pricing mechanisms.

In summary, there is support for the principal-agent hypotheses linking contract form with (a) information systems (Conlon & Parks, 1988; Ec-cles, 1985; Eisenhardt, 1985), (b) outcome uncer-tainty (Eisenhardt, 1985), (c) outcome measur-ability (Anderson, 1985; Eisenhardt, 1985), (d) time (Conlon & Parks, 1988), and (e) task pro-grammability (Eccles, 1985; Eisenhardt, 1985). Moreover, this support rests on research using a variety of methods including questionnaires, secondary sources, laboratory experiments, and interviews.

## Recommendations for Agency Theory Research

As argued above, agency theory makes con-tributions to organization theory, is testable, and has empirical support. Overall, it seems reason-able to urge the adoption of an agency theory perspective when investigating the many prob-lems that have a principal-agent structure. Five specific recommendations are outlined below for using agency theory in organizational re-search.

### Focus on Information Systems, Outcome Uncertainty, and Risk

McGrath, Martin, and Kukla (1981) argued that research is a knowledge accrual process. Using this accrual criterion, next steps for agency theory research are clear: *Researchers should focus on information systems, outcome uncertainty, and risk.* These agency variables make the most unique contribution to organiza-tional research, yet they have received little em-pirical attention (Table 3). It is important that re-searchers place emphasis on these variables in order to advance agency theory and to provide new concepts in the study of familiar topics such as impression management, innovation, verti-cal integration, compensation, strategic alli-ances, and board relationships.

Studying risk and outcome uncertainty is par-ticularly opportune because of recent advances in measuring risk preferences. By relying on the works of Kahneman and Tversky (1979), Mac-Crimmon and Wehrung (1986), and March and Shapira (1987), the organizational researcher can measure risk preference more easily and realistically. These techniques include direct measures of risk preference such as lotteries and indirect measures using demographic charac-teristics such as age and wealth and payoff characteristics such as gain versus loss. (See March and Shapira, 1987, for a review.)

### Key on Theory-Relevant Contexts

Organizational theory usually is explored in settings in which the theory appears to have greatest relevance. For example, institutional and resource dependence theories were devel-oped primarily in large, public bureaucracies in which efficiency may not have been a pressing

70

concern. The recommendation here is to take the same approach with agency theory: *Key on theory-relevant contexts*.

Agency theory is most relevant in situations in which contracting problems are difficult. These include situations in which there is (a) substantial goal conflict between principals and agents, such that agent opportunism is likely (e.g., owners and managers, managers and professionals, suppliers and buyers); (b) sufficient outcome uncertainty to trigger the risk implications of the theory (e.g., new product innovation, young and small firms, recently deregulated industries); and (c) unprogrammed or team-oriented jobs in which evaluation of behaviors is difficult. By emphasizing these contexts, researchers can use agency theory where it can provide the most leverage and where it can be most rigorously tested. Topics such as innovation and settings such as technology-based firms are particularly attractive because they combine goal conflict between professionals and managers, risk, and jobs in which performance evaluation is difficult.

**Expand to Richer Contexts**

Perrow (1986) and others have criticized agency theory for being excessively narrow and having few testable implications. Although these criticisms may be extreme, they do suggest that research should be undertaken in new areas. Thus, the recommendation is *to expand to a richer and more complex range of contexts*.

Two areas are particularly appropriate. One is to apply the agency structure to organizational behavior topics that relate to information asymmetry (or deception) in cooperative situations. Examples of such topics are impression management (Gardner & Martinko, 1988), lying and other forms of secrecy (Sitkin, 1987), and blame (Leatherwood & Conlon, 1987). Agency theory might contribute an overall framework in which to place these various forms of self-interest, leading to a better understanding of when such behaviors will be likely and when they will be effective.

The second area is expansion beyond the pure forms of behavior and outcome contracts as described in this article to a broader range of contract alternatives. Most research (e.g., Anderson, 1985; Eisenhardt, 1985, 1988) treats contracts as a dichotomy: behavior versus outcome. However, contracts can vary on a continuum between behavior and outcome contracts. Also, current research focuses on a single reward, neglecting many situations in which there are multiple rewards, differing by time frame and contract basis. For example, upper level managers usually are compensated through multiple rewards such as promotions, stock options, and salary. Both multiple and mixed rewards (behavior and outcome) present empirical difficulties, but they also mirror real life. The richness and complexity of agency theory would be enhanced if researchers would consider this broader spectrum of possible contracts.

**Use Multiple Theories**

A recent article by Hirsch et al. (1987) eloquently compared economics with sociology. They argued that economics is dominated by a single paradigm, price theory, and a single view of human nature, self-interest. In contrast, the authors maintained that a strength of organizational research is its polyglot of theories that yields a more realistic view of organizations.

Consistent with the Hirsch et al. arguments, the recommendation here is *to use agency theory with complementary theories*. Agency theory presents a partial view of the world that, although it is valid, also ignores a good bit of the complexity of organizations. Additional perspectives can help to capture the greater complexity.

This point is demonstrated by many of the empirical studies reviewed above. For example, the Singh and Harianto (in press) and Kosnik (1987) studies support agency theory hypotheses, but they also use the complementary perspectives of hegemony and managerialism. These perspectives emphasize the power and political aspects of golden parachutes and green-

71

mail, respectively. Similarly, the studies by Eisenhardt (1988) and Conlon and Parks (1988) combine institutional and agency theories. The institutional emphasis on tradition complements the efficiency emphasis of agency theory, and the result is a better understanding of compensation. Other examples include Anderson (1985), who coupled agency and transaction cost, and Eccles (1985), who combined agency with equity theory.

### Look Beyond Economics

The final recommendation is *that organizational researchers should look beyond the economics literature*. The advantages of economics are careful development of assumptions and logical propositions (Hirsch et al., 1987). However, much of this careful theoretical development has already been accomplished for agency theory. For organizational researchers, the payoff now is in empirical research, where organizational researchers have comparative advantage (Hirsch et al., 1987). To rely too heavily on economics with its restrictive assumptions such as efficient markets and its single-perspective style is to risk doing second-rate economics without contributing first-rate organizational research. Therefore, although it is appropriate to monitor developments in economics, it is more useful to treat economics as an adjunct to more mainstream empirical work by organizational scholars.

## Conclusion

This paper began with two extreme positions on agency theory—one arguing that agency theory is revolutionary and a powerful foundation (Jensen, 1983) and the other arguing that the theory addresses no clear problem, is narrow, lacks testable implications, and is dangerous (Perrow, 1986). A more valid perspective lies in the middle. Agency theory provides a unique, realistic, and empirically testable perspective on problems of cooperative effort. The intent of this paper is to clarify some of the confusion surrounding agency theory and to lead organizational scholars to use agency theory in their study of the broad range of principal-agent issues facing firms.

## References

Anderson, E. (1985) The salesperson as outside agent of employee: A transaction cost analysis. *Marketing Science*, 4, 234–254.

Amihud, Y., & Lev, B. (1981) Risk reduction as a managerial motive for conglomerate mergers. *Bell Journal of Economics*, 12, 605–616.

Argawal, A., & Mandelker, G. (1987) Managerial incentives and corporate investment and financing decisions. *Journal of Finance*, 42, 823–837.

Arrow, K. (1971) *Essays in the theory of risk bearing*. Chicago: Markham.

Barnard, C. (1938) *The functions of the executive*. Cambridge, MA: Harvard University Press.

Barney, J. (1988) *Agency theory, employee stock ownership and a firm's cost of equity capital*. Unpublished working paper, Texas A&M University, College Station.

Barney, J., & Ouchi, W. (Eds.) (1986) *Organizational economics*. San Francisco: Jossey-Bass.

Basu, A., Lal, R., Srinivasan, V., & Staelin, R. (1985) Salesforce compensation plans: An agency theoretic perspective. *Marketing Science*, 4, 267–291.

Berle, A., & Means, G. (1932) *The modern corporation and private property*. New York: Macmillan.

Bolton, M. (1988) *Organizational miming: When do late adopters of organizational innovations outperform pioneers?* Paper presented at the meeting of the Academy of Management, Anaheim, CA.

Burt, R. (1979) A structural theory of interlocking corporate directorates. *Social Networks*, 1, 415–435.

Capitalism in the making. (1984, April 30) *Time*, p. 62.

Chandler, A. (1962) *Strategy and structure*. New York: Doubleday.

Conlon, E., & Parks, J. (1988) The effects of monitoring and tradition on compensation arrangements: An experiment on principal/agent dyads. In F. Hoy (Ed.), *Best papers pro-*

72

*ceedings* (pp. 191–195). Anaheim, CA: Academy of Management.

Cyert, R., & March, J. (1963) *A behavioral theory of the firm.* Englewood Cliffs, NJ: Prentice-Hall.

DeJong, D., Forsythe, R., & Uecker, W. (1985) Ripoffs, lemons and reputation formation in agency relationships: A laboratory market study. *Journal of Finance,* 50, 809–820.

Demski, J. (1980) *A simple case of indeterminate financial reporting.* Working paper, Stanford University.

Demski, J., & Feltham, G. (1978) Economic incentives in budgetary control systems. *Accounting Review,* 53, 336–359.

Dornbusch, S., & Scott, W. R. (1974) *Evaluation and the exercise of authority.* San Francisco: Jossey-Bass.

Eccles, R. (1985) Transfer pricing as a problem of agency. In J. Pratt & R. Zeckhauser (Eds.), *Principals and agents: The structure of business* (pp. 151–186). Boston: Harvard Business School Press.

Eisenhardt, K. (1985) Control: Organizational and economic approaches. *Management Science,* 31, 134–149.

Eisenhardt, K. (1988) Agency and institutional explanations of compensation in retail sales. *Academy of Management Journal,* 31, 488–511.

Fama, E. (1980) Agency problems and the theory of the firm. *Journal of Political Economy,* 88, 288–307.

Fama, E., & Jensen, M. (1983) Separation of ownership and control. *Journal of Law and Economics,* 26, 301–325.

Galbraith, J. (1973) *Designing complex organizations.* Reading, MA: Addison-Wesley.

Gardner, W., & Martinko, M. (1988) Impression management: An observational study linking audience characteristics with verbal self-presentations. *Academy of Management Journal,* 31, 42–65.

Gausch, J., & Weiss, A. (1981) Self-selection in the labor market. *American Economic Review,* 71, 275–284.

Harris, M., & Raviv, A. (1978) Some results on incentive contracts with application to education and employment, health insurance, and law enforcement. *American Economic Review,* 68, 20–30.

Harris, M., & Raviv, A. (1979) Optimal incentive contracts with imperfect information. *Journal of Economic Theory,* 20, 231–259.

Hirsch, P., & Friedman, R. (1986) Collaboration or paradigm shift? Economic vs. behavioral thinking about policy? In J. Pearce & R. Robinson (Eds.), *Best papers proceedings* (pp. 31–35). Chicago: Academy of Management.

Hirsch, P., Michaels, S., & Friedman, R. (1987) "Dirty hands" versus "clean models": Is sociology in danger of being seduced by economics? *Theory and Society,* 317–336.

Holmstrom, B. (1979) Moral hazard and observability. *Bell Journal of Economics,* 10, 74–91.

Jensen, M. (1983) Organization theory and methodology. *Accounting Review,* 56, 319–338.

Jensen, M. (1984) Takeovers: Folklore and science. *Harvard Business Review,* 62(6), 109–121.

Jensen, M., & Meckling, W. (1976) Theory of the firm: Managerial behavior, agency costs, and ownership structure. *Journal of Financial Economics,* 3, 305–360.

Jensen, M., & Roeback, R. (1983) The market for corporate control: Empirical evidence. *Journal of Financial Economics,* 11, 5–50.

Kahneman, D., & Tversky, A. (1979) Prospect theory: An analysis of decisions under risk. *Econometrica,* 47, 263–291.

Kosnik, R. (1987) Greenmail: A study in board performance in corporate governance. *Administrative Science Quarterly,* 32, 163–185.

Lambert, R. (1983) Long-term contracts and moral hazard. *Bell Journal of Economics,* 14, 441–452.

Lawrence, P., & Lorsch, J. (1967) *Organization and environment.* Boston: Division of Research, Harvard Business School.

Leatherwood, M., & Conlon, E. (1987) Diffusibility of blame: Effects on persistence in a project. *Academy of Management Journal,* 30, 836–848.

MacCrimmon, K., & Wehrung, D. (1986) *Taking risks: The management of uncertainty.* New York: Free Press.

March, J. (1962) The business firm as a political coalition. *Journal of Politics,* 24, 662–678.

March, J., & Shapira, Z. (1987) Managerial perspectives on risk and risk taking. *Management Science,* 33, 1404–1418.

March, J., & Simon, H. (1958) *Organizations.* New York: Wiley.

McGrath, J., Martin, J., & Kukla, R. (1982) *Judgment calls in research.* Beverly Hills, CA: Sage.

Meet Mike Jensen, the professor of merger mania. (1988, February 8) *Business Week,* pp. 66–67.

Mitnick, B. (1986) *The theory of agency and organizational analysis.* Unpublished working paper, University of Pittsburgh.

Ouchi, W. (1979) A conceptual framework for the design of organizational control mechanisms. *Management Science,* 25, 833–848.

Pearce, J., Stevenson, W., & Perry, J. (1985) Managerial compensation based on organizational performance: A time series analysis of the effects of merit pay. *Academy of Management Journal,* 28, 261–278.

73

198

Perrow, C. (1986) *Complex organizations*. New York: Random House.

Pettigrew, A. (1973) *The politics of organizational decision making*. London: Tavistock.

Pfeffer, J. (1981) *Power in organizations*. Marshfield, MA: Pittman.

Pfeffer, J., & Salancik, G. (1974) Organizational decision making as a political process: The case of a university budget. *Administrative Science Quarterly*, 19, 135–151.

Ross, S. (1973) The economic theory of agency: The principal's problem. *American Economic Review*, 63, 134–139.

Shavell, S. (1979) Risk sharing and incentives in the principal and agent relationship. *Bell Journal of Economics*, 10, 53–73.

Singh, H., & Harianto, F. (in press) Management-board relationships, takeover risk and the adoption of golden parachutes: An empirical investigation. *Academy of Management Journal*.

Sitkin, S. (1987) *Secrecy in organizations: The limits of legitimate information control*. Working paper, University of Texas, Austin.

Spence, A. M. (1975) The economics of internal organization: An introduction. *Bell Journal of Economics*, 6, 163–172.

Spence, A. M., & Zeckhauser, R. (1971) Insurance, information, and individual action. *American Economic Review*, 61, 380–387.

Thompson, J. (1967) *Organizations in action*. New York: McGraw-Hill.

Ungson, G., & Steers, R. (1984) Motivation and politics in executive compensation. *Academy of Management Review*, 9, 313–323.

Walker, G., & Weber, D. (1984) A transaction cost approach to make-or-buy decisions. *Administrative Science Quarterly*, 29, 373–391.

Walking, R., & Long, M. (1984) Agency theory, managerial welfare, and takeover bid resistance. *The Rand Journal of Economics*, 15, 54–68.

White, H. (1985) Agency as control. In J. Pratt & R. Zeckhauser (Eds.), *Principals and agents: The structure of business* (pp. 187–214). Boston: Harvard Business School Press.

Williamson, O. (1975) *Markets and hierarchies: Analysis and antitrust implications*. New York: Free Press.

Wilson, R. (1968) On the theory of syndicates. *Econometrica*, 36, 119–132.

Wolfson, M. (1985) Empirical evidence of incentive problems and their mitigation in oil and gas shelter programs. In J. Pratt & R. Zeckhauser (Eds.), *Principals and agents: The structure of business* (pp. 101–126). Boston: Harvard Business School Press.

Zenger, T. (1988) *Agency sorting, agent solutions and diseconomies of scale: An empirical investigation of employment contracts in high technology R&D*. Paper presented at the meeting of the Academy of Management, Anaheim, CA.

*Kathleen M. Eisenhardt (Ph.D., Stanford University) is Assistant Professor at Stanford University. Correspondence can be sent to her at the Department of Industrial Engineering and Engineering Management, 346 Terman Building, Stanford University, Stanford, CA 94305.*

74

# AN AGENCY THEORY VIEW OF THE MANAGEMENT
# OF END-USER COMPUTING

**Vijay Gurbaxani**
Graduate School of Management
University of California, Irvine

**Chris F. Kemerer**
Sloan School of Management
Massachusetts Institute of Technology

## ABSTRACT

The growth in end-user computing (EUC) in organizations and its implications for the degree of centralization of the information services function have led to the need for a theory that will assist in the management of this process. This paper employs microeconomics and, in particular, agency theory to describe the development of EUC in organizations. The results suggest that agency theory provides useful insights and significant normative implications for the management of computing in organizations.

## 1. INTRODUCTION

The dramatic decline in the costs of hardware and the trend towards the increased power of microcomputers and minicomputers have enabled significant growth in end-user computing (EUC). This trend has implications not only for the management of EUC but also for the degree of centralization of the Information Systems (IS) function in organizations. Therefore, there has been increased focus on the *organizational* issues surrounding EUC, as evidenced by senior IS executives' responses in several recent surveys.[1] Management issues related to decentralization of the IS department also ranked high on their list of concerns.

This interest in EUC has resulted in numerous articles in the academic and practitioner literature. The primary thrust of many of these articles is prescriptive and suggests alternative managerial strategies for EUC (Alavi, Nelson and Weiss 1987; Gerrity and Rockart 1986; Henderson and Treacy 1985; Munro, Huff and Moore 1987). Some studies have analyzed the characteristics of end-users and their tasks[2] and how these tasks evolve (Huff, Munro and Martin 1988). Robey and Zmud (1989) have recently criticized the EUC literature for not being "grounded in specific theories of organizational behavior." The current paper proposes the use of *agency theory* as a theoretical base and integrative approach within which to understand the EUC phenomenon.

The definition of EUC adopted in this paper is that of Davis and Olson (1985), namely, "the capability of users to have direct control of their own computing needs." This definition of EUC emphasizes the control aspects of the problem which, it will be argued later, are at the heart of the issue. In particular, it highlights the division of control between the end-user departments and the central IS organization.

The reference discipline employed in this paper is microeconomics encompassing agency theory, as originally suggested by Kriebel and Moore (1980). While traditional microeconomics has proven useful in analyzing a large variety of problems, it has not been widely used in analyzing intra-firm managerial control problems due to its assumptions of costless information transfer and of goal congruence of managers within the firm. Agency theory extends the microeconomic approach by relaxing these assumptions and, therefore, will be shown to be particularly appropriate for the intra-firm nature of the EUC control problem. The theory developed here has significant normative implications for the management of computing in organizations.

The outline of this paper is as follows. The research problem and approach are presented in Section 2. Section 3 introduces the principal-agent problem in IS within a microeconomic framework and analyzes its impact on the production strategies for information services. Managerial implications and concluding remarks are presented in Section 4.

## 2. THE RESEARCH PROBLEM AND APPROACH

This section begins with an introduction of the research problem – control of the provision of IS services. Next, the salient features of the traditional microeconomic approach and its shortcomings in analyzing managerial behavior in this context are presented. This is followed by a brief discussion of agency theory, which extends the traditional microeconomic approach to address these deficiencies.

## 2.1 Research Problem

Given the nature of the supply of and demand for information services, the organization must determine how the internal provision of information services will be organized so as to maximize the net value of information services. The focus here is on the *control* issues related to the internal provision of these services. The definition of control adopted in this paper is that of Fama and Jensen (1983), namely, "the ability to i) choose the decision initiative to be implemented and ii) to measure the performance of agents and implement a reward structure." Control issues that govern IS activities include the choice of the organization structure of the IS department, managerial compensation contracts, the decision to mandate that a service be acquired from the central IS group, chargeback systems for information services, and budget allocation mechanisms. The choice of control mechanisms is naturally a major determinant of the effectiveness of IS activities.

The authority to determine how a specific activity is performed is termed here as a *decision right*. Formal modeling approaches recognize that initially all decision rights reside with top management, who may decide to allocate some or all of these rights to IS and end-user departments. Then, the locus of control is determined by the partitioning of decision rights between the different members of the organization. Thus, the control problem may be viewed as determining the optimal partition of these decision rights.

Decentralized computing, defined here as the transfer of control from centralized IS departments to end-user or functional departments, has continued to grow in scope and in degree (Arthur Andersen 1986). The decentralization of computing cannot be explained simply by examining the economics of the production of information services. If that were true, one might witness the growth of *distributed computing* as distinguished from *decentralized computing*. Distributed computing is defined as the location of hardware, software and personnel at various sites throughout the organization with the important provision that control decision rights remain vested in a central authority.

An underlying factor in the growth of decentralized computing was the dissatisfaction of users with the centralized environment. In theory, it is feasible to develop a centralized plan for the provision of information services wherein all users are satisfied. Yet, this has rarely occurred. It shall be argued below that principal-agent problems have been a significant factor in decreasing the likelihood of success of a centralized IS approach. However, before developing this argument, the traditional microeconomics argument is presented to provide a framework with which to build the agency model.

## 2.2 The Traditional Microeconomics View

The microeconomics approach to developing positive or descriptive models of a phenomenon assumes net value maximizing behavior. To develop a positive theory of IS management, one would build a model of this process by assuming that practices relating to the management of computing are an outcome of net value maximizing behavior (Silberberg 1978). Thus, one would assume that the goal of the firm would be to maximize the net value of information services to the organization and derive, for example, the testable implication that in the early years of computing, firms centralized computing services to exploit economies of scale in hardware. These hypotheses of managerial behavior could then be tested against empirical data to determine the validity of the model.

One traditional microeconomic approach has been to assume a number of ideal conditions, under which it has been shown that optimizing behavior on the parts of individuals and firms under pure-competition leads to a Pareto-optimal social outcome, i.e., one where no other allocation makes all parties concerned at least as well off and one or more parties better off (Hirshleifer 1980). It implies that, under certain conditions, social welfare is maximized simply as a result of the individual economic players acting out of self-interest. More formally, a Pareto-optimal allocation results in a competitive equilibrium implying efficiency among consumers in the allocation of consumption goods, efficiency among resource owners in the provision of resources for productive uses, and efficiency among firms in the conversion of resources into consumable goods.

While the above discussion applies to economic actors in a competitive market, it can also be extended to apply within a firm. In the context of the management of IS, the parallel situation would be the creation of a market for information services within a firm (perhaps even including economic actors outside the firm). Thus, one could consider a situation where individual departments would be allowed to act as consumers or suppliers of information services. If the net value of information services to the organization were maximized using such an approach, the task facing the firm would be the creation and maintenance of such a market.

However, several factors may cause a market failure where social welfare is not maximized in a market situation. These include the presence of market power and the existence of externalities. Market power is usually seen as monopoly or monopsony power. Externalities occur when the actions of an economic agent affect the interests of other agents in a way not captured by market prices. Both of these factors limit the applicability of the traditional microeconomics results and may lead to situations where a pure market-based approach is inadequate. Vertical integration is often cited as a possible solution to these problems, since it allows the internalization of externalities

and limits market power. In this paper, it is argued that i) both market power and externalities are present in the intra-firm IS context and that ii) market power is exercised and actions that cause externalities are taken because of problems due to the *agency relationships* (discussed below) among the actors within a firm.

## 2.3 The Theory of Agency

An agency relationship can be said to occur whenever one party depends on the actions of another party. More formally, Jensen and Meckling (1976) define an agency relationship as "a contract under which one or more persons (the principal(s)) engage another person (the agent) to perform some service on their behalf which involves delegating some decision making authority to the agent." In an organizational context, a firm hires employees (agents) in part to exploit economies of specialization. Yet, these employees often act in a manner that is inconsistent with maximizing the welfare of the firm. Agency theory argues that this occurs because

(a) the goals of the principal and the agent are often inconsistent with one another ("goal incongruence") and

(b) the principal cannot perfectly and costlessly monitor the *actions* and the *information* of the agent ("information asymmetries").

Since agents are usually better informed than their principals about their tasks, organizations would do better if all information could be shared at zero cost, or if there was no divergence between the goals of the principals and the agents. The economic loss that occurs due to the absence of such optimal conditions is called the *agency cost*. The components of agency costs are *monitoring costs* expended by the principal to observe the agent, *bonding costs* incurred by the agent to make his or her services more attractive, and *residual loss*, which are the opportunity costs borne by the principal due to the difference in outcomes that would obtain between the principal's and agent's execution of the task (Jensen and Meckling 1976). An implication of the assumption of net value maximization and the existence of agency costs is that the principal seeks to minimize agency costs through the use of control mechanisms. The primary control mechanisms in organizations are the performance measurement and evaluation system, the reward and punishment system, and the system for assigning decision rights among participants in the organization (Jensen 1983).

In the case of costless information transfer and the absence of agency costs, as is assumed by the traditional microeconomics approach, the control problem is inconsequential. One can simply assume that all information that a central planner requires to make a decision and that is possessed by other actors within the firm can be acquired without

cost. Further, since all actors behave in a manner that is consistent with maximizing the value of the firm, no control mechanisms are required to ensure the consistency of managerial behavior with the goals of the firm. However, in a realistic setting, the control problem assumes importance because of the existence of information asymmetries and goal incongruencies and the resulting agency costs.

Eisenhardt (1989) has articulated the usefulness of agency theory in analyzing managerial problems characterized by goal conflicts, outcome uncertainty, and unprogrammed or team-oriented tasks. Many IS activities fit this description, and it has been suggested that a large number of organizational problems in the management of IS can be analyzed successfully in an agency context (Gurbaxani and Kemerer 1989; Beath and Straub 1989; Robey and Zmud 1989; Klepper 1990). The design of effective control mechanisms for IS activities is particularly difficult, since the agency relationship occurs in a dynamic, rapidly changing environment and management practices have little time to stabilize (Nolan 1979; Gurbaxani and Mendelson 1990). In this paper, the focus is on the impact of agency costs on the organization of the *internal* provision of information services.

An alternative approach would be transaction cost economics, an approach with similarities to agency theory in its emphasis on information and uncertainty (Williamson 1985). However, as noted by Eisenhardt (1989), agency theory distinguishes itself from transaction cost theory by its inclusion of the notions of risk aversion and information as a commodity.

## 3. AN AGENCY VIEW OF INFORMATION SERVICES

The key issues that arise in an agent-theoretic analysis of the management of IS are an identification of the economic actors and their objectives, an analysis of how these objectives result in conflict, and an analysis of the nature of the resulting organizational costs. These issues must be considered in conjunction with the microeconomic and technological characteristics of the IS environment to determine the optimal strategies for the management of IS resources. Specifically examined are the impact of agency costs on the growth of EUC and the implications for the degree of centralization of the information services functions.

### 3.1 The Agency Structure of Traditional Computing in Organizations

In order to provide a model of current end-user computing, it is helpful to begin with a brief discussion of traditional computing in organizations to show the origins of EUC. The level of analysis is the department and three types of economic units will be relevant: top management,

the centralized IS department, and end-user departments (see Figure 1).[3]



**Figure 1. Agency Relationships**

There are three resulting principal-agent relationships. In two of these relationships, the principal is top management and the agents are the functional departments and the IS department. In the third relationship, each end-user department is a principal and the IS department is the agent.[4] The objectives of each of these actors are considered in turn, focusing on the IS aspects of the principal-agent relationships. It will be argued in Section 3.2 that the individual objectives of each of these actors can be in conflict with one another and result in agency costs. However, before coming to that conclusion, it is useful to examine how this structure for providing information services within the organization came about.

When computing was first introduced into organizations, most end-users and top management, specifically, were unfamiliar with the technology. This resulted in top management creating IS departments and hiring specialists in the production of information services. For the same reason, most decision rights related to the management of computing were allocated to the IS department. The decision to centralize computing was driven primarily by the costs of computing.[5] The demand generated by any single end-user group often did not justify such a large investment. Thus, the demands of various end-user groups had to be aggregated to justify the investment. The decision rights related to hardware and software selection were typically located in the IS department. Applications software was developed almost exclusively by professionals who were located in the IS department. Since individual end-user departments were uncertain of their future demands for software services, the appropriate strategy for the location of software professionals was to centralize the programming function since this would simplify the management of these professionals.

Therefore, the centralization of computing was a result of organizations seeking to exploit economies of scale and of specialization that were warranted by the high costs of computing. Due to supply-side considerations in that time, the costs of production outweighed any other costs in determining the strategy for the provision of information services. The problem associated with this shared resource approach is that the socially optimum solution may not be any user group's local optimum. This idea is critical to the discussion below of the impacts of agency costs on the provision of information services.

### 3.2 Market Failures in Organizational Computing Due to Agency

As the unit costs of computing decreased over time, and as minicomputer and microcomputer technology became available, decentralized computing became feasible, as will be seen below. *However, the changing economics of information systems supply are a necessary but not sufficient condition for decentralized computing, as opposed to merely distributed computing.* To see this, the nature of agency costs in a centralized environment are discussed below.

#### 3.2.1 IS Department as an Agent of Top Management

In the traditional environment, top management relied upon IS specialists as their agents to provide IS services. These agents were typically organized into one centralized department due to the economies of scale and specialization noted above. However, this agency relationship introduces costs to the organization through goal incongruencies and information asymmetries.

While net value maximization of information services may be the desired intent of top management for IS managers, the IS managers' actual behavior patterns sometimes suggest that their "objective function" may be quite different. For example, the salaries of these managers are often related to the scale of the operation, inducing them to indulge in so-called "empire building." A related cost arises because of the value managers place on the control of a resource that may increase their political power within the organization. Another problem is termed the "asymmetric cost" problem (Mendelson 1990). Here, managers often make sub-optimal decisions because the cost of the decision to the manager may be quite different than that incurred by the firm. For example, a manager's evaluation is sometimes based on the quality of services provided rather than on its cost effectiveness. This is often stated in the practitioner literature as "No one ever got fired for buying IBM." This is an example of the risk-averse nature of the IS manager-agent. IS managers also often suffer from the "professional syndrome" (Mendelson 1990), wherein they have incentives to acquire the newest hardware and software technologies with insufficient regard for cost justification. This is consistent with maximizing

behavior of the IS professional whose market value is partly determined by his familiarity with new technologies.

The optimal allocation of information services typically requires that the marginal value of information services to a division equal the marginal cost of providing these services (Hirshleifer 1980). If information transfer were costless, one could assume that the IS department possessed both the cost and value information required to implement such an allocation. Thus, information asymmetries would not be an issue. Furthermore, since their actions would be completely known to top management, the IS department could be expected to maximize the net value of information services to the firm.

However, the existence of asymmetric information precludes such a solution. The primary information asymmetry in the IS context is that knowledge of the value of a given IS task is almost always possessed by the end-user, while information about the execution of the task is possessed by the IS department. This information asymmetry also extends to top management, who are neither completely aware of the value of information generated by IS activities to the end-user departments nor of the cost and technological information possessed by the IS department. Thus, top management is faced with the problem of constructing a control system that will maximize the net value of information services to the firm while taking into account the existence of these information asymmetries.

Top management traditionally imposed one of two control structures: a profit center approach or a cost center approach. In a profit center, the performance of the IS manager is measured by the magnitude of profits that he or she generates, while in the case of a cost center, the performance metrics are related to adherence to budgets or by comparison with "standard costs." Each of these creates very different sets of incentives for the IS manager. The profit center encourages efficiency in the production of information services, but also creates incentives for the IS manager to act as a monopolist to increase profits. This, in turn, raises the likelihood that the prices of computing services will be higher than optimal. The cost center, on the other hand, does not create incentives for higher prices, but neither does it encourage efficient production.[6]

In both of these control structures, the welfare of the organization is reduced by the agency costs that result from the actions of the IS manager. When the IS department is set up as a cost center, the costs result from the inefficient production of information services. These costs are manifested as delays in operations, backlogs in software development, and higher total costs (as distinguished from unit costs) for information services. In the case of a profit center, such costs are incurred primarily as higher monopoly prices rather than as free-market prices for services.

As the costs of computing continued to decrease over time, and as minicomputer technology became a feasible option, the decision not to mandate that all computing services be acquired from central IS meant that individual end-user departments were given the right to implement decentralized computing. The fact that decentralized computing was implemented by some end-users – even though it was initially more expensive than centralized computing services due to the fixed costs and lost economies of scale and specialization – strongly suggests that these end-users were incurring costs beyond those seen in accounting statements. This suggests that end-user departments may have exercised this option in part to minimize the agency costs resulting from the self-interested behavior of the IS department. Decentralized computing can be seen as an effective means of limiting the market power of IS departments.

### 3.2.2 User Department as an Agent of Top Management

Analogous to the IS department's role as an agent to the firm, each end-user or functional department also acts as an agent (Figure 1). Therefore, their behavior also reflects goal incongruencies and information asymmetries in their relationship with the top management principal. The discussion here, however, will be limited to the effect of these factors on the allocation of IS resources within the firm.

Decisions that maximize the net value of information services to the firm may not be locally optimal, that is, they might not maximize the net value of information services to the individual end-user departments. End-users may be dissatisfied with resource allocation decisions. For example, in a mainframe acquisition decision where there are many possible end-users, each set of end-users may prefer a different type of machine. In the case where there is insufficient demand to justify the purchase of more than one machine, only a subset of end-users will receive the machine of first choice, and others will have to make do with a lower ranking choice.[7] Similar situations arise in the acquisition of software packages as well.

The analogous situation exists in the case of a software development task. The globally optimal specifications for such a task may be an outcome of meeting the demands of numerous end-user groups. Individual end-user groups would prefer customized applications that, in all likelihood, would also have better performance, since they would not be constrained by the requirements of other end-users. Moreover, end-users whose application development requests are queued behind others of higher value to the organization incur waiting costs.

End-user departments also have incentives that encourage them to control their own information. There are several possible reasons for this. The possession of information

that is of significant value to the firm often results in increased power to the owner of the information. Another reason may be that the information may allow top management to monitor the performance of an end-user department more closely, a possibly undesirable situation for the end-user.

In all of the above situations, undoubtedly some end-user departments could be made better off if the resource allocation decisions were modified in their favor. Therefore, the end-users now perceive that they can increase their welfare by biasing the information they provide the IS department to increase the likelihood of a more favorable outcome. For example, an end-user may request a higher priority on a timesharing machine than is really warranted by the task or may demand a more powerful personal computer than the one that is the most cost-effective. In such cases, the cost imposed on other end-users stems from a reduction in resources available to them. Given the assumption of self-interested behavior, such costs are likely.

Of course, end-users are subject to monitoring by top management that limits the amount of bias in information that they can provide. However, monitoring is rarely perfect, and engaging in monitoring activities also results in monitoring costs to the organization. The net result is that resource allocation schemes that are in some part dependent on the full disclosure of information by agents are unlikely to be totally successful in practice. The challenge facing the organization is to develop a control strategy that aligns the self-interest of agents with the interests of the firm. Agency theory suggests mechanisms, known as incentive compatible contracts, for managing such problems, and examples of such approaches will be discussed in Section 4.

### 3.2.3 IS Department as an Agent of End-User Departments

The third and final agency relationship, consistent with the IS department being a "staff" as opposed to "line" function in most organizations, is that of the centralized IS department acting as an agent for an end-user department. This relationship also provides for a strong additional source of conflict within the organization. The IS department is effectively the agent of multiple principals (i.e., the top management principal and the end-user principal) whose goals may not converge, as has already been discussed in Section 3.2.2. Added to this may be the IS department's own agenda (the potential conflict with top management having been discussed in Section 3.2.1). Therefore, conflict between the IS department and an end-user department can come about because a) the IS department is trying to act as an agent for top management and, therefore, may not act in accordance with the desired behavior of the self-interested end-user and/or b) the IS department is itself

engaging in self-interested behavior at the expense of the user department.

There are several forms that these goal incongruencies may take in the IS context. Assuming that the IS department is acting on behalf of the organization, then they will be providing software systems and hardware services that meet the needs of the entire organization, not just an individual end-user department. Therefore, a decision that IS may make on behalf of the organization may be suboptimal for any given department. In particular, the IS department will engage in activities to promote the long-term computing environment serving a variety of end-users. Therefore, any particular end-user will bear additional costs, including delay costs and integration costs, because they are using shared resources.

For example, consider the issue of integration in software development. As the goal of central IS is to support the needs of the entire organization, the need for integration is clear and vital. Also, as a central provider of services, IS can exploit scale economies by developing policies and procedures that provide a consistent and integrated base, such as a central database, development platforms or interface standards. End-user departments may have neither the incentive nor the scale to justify this type of effort. Moreover, with a centralized control mechanism, redundant efforts are less likely to occur, a result that is consistent with the goals of the organization.

This divergence of goals has been noted by several EUC researchers in terms of the lack of effort expended toward integration and coordination. For example, in separate studies both Guimaraes (1984, p. 5) and Alavi (1985, p. 17) have noted that an end-user over-emphasis on short-term operational issues at the expense of longer-term managerial concerns has led to many EUC problems with lack of systems integration.

Another aspect to this shared resource phenomenon is that it may appear to be a public good to the end-user. Given self-interested behavior on the part of the end-user, it is expected that they will tend to use more of the IS resource than might be organizationally desired if the control mechanisms do not insure that the end-user fully bears the costs of such consumption.

Of course, the IS department may not always act in accordance with what the top management-principal may desire either. Given the difficulty in assessing the value of IS services, many organizations may treat it as a "utility," where the IS department management is evaluated essentially on the ability to deliver a consistent quality of service. This might be implemented by metrics such as machine or network uptime, or low levels of end-user problem reports. In this type of environment, IS department management can become very risk averse, as changes may involve disruptions in service levels. Therefore, any end-user's desire for applications or technologies that differ

from past approaches may be discouraged. This phenomenon is particularly relevant in IS services due to the rapid rate of technological change in this area.

Consider the following example from applications development illustrating the issue of risk aversion on the part of the central IS-agent. Traditionally, large systems have been developed using the systems development life cycle (SDLC), a process designed to initially elicit system requirements from end-users, and then to build systems in a carefully planned series of sequential stages that emphasize system validity, correctness and maintainability, rather than speed of development. An alternative approach is prototyping, which allows shorter lead times for the delivery of a limited set of functionality. In prototyping, development work continues until the user is satisfied. Thus, prototyping is essentially an outcome-based control strategy, while the systems development life cycle, with its extensive task checklists, is essentially an input, or behavior-based approach (Ouchi 1979; Eisenhardt 1985). Agency theory would predict that the risk-averse agent (central IS) would prefer a behavior-based approach, since the outcome-based approach entails greater risk. Conversely, the end-user principal, who cannot perfectly monitor the agent's behavior, would prefer an outcome-based approach. In fact, these preferences are observed in practice. For example, Rockart and Flannery (1986, p. 288), in their study of EUC, note that end-users find central IS's tools, methods and processes "entirely inappropriate" for a significant part of their new applications.

In summary, conflict between the end-user principal and the IS department-agent can develop from either the IS department role in representing its top management principal or due to the goals of the IS department itself.

## 3.3 Conclusions

Given the above discussion, a model of the provision of information services must incorporate the behavioral assumptions that the goals of principals and agents may diverge and that agents act out of self-interest. It should also recognize that information transfer is costly, and moreover, it cannot be presumed that an agent will be willing to reveal private information if such revelation is inconsistent with his or her goals.

Based on the above agency model of IS provision, there exists an essential tension between the centralization and decentralization of IS services. Existing centralized IS departments will prefer the status quo, in part to maximize their own welfare. End-users will desire greater autonomy over their computing, in part in order to avoid the externalities and agency costs which arise in the centralized solution. Into this environment comes the technical feasibility of end-user computing. This provides an option for end-users to provide at least some of their own IS services. That this option has been acted upon in practice suggests support for the agency model.

## 4. MANAGERIAL IMPLICATIONS

### 4.1 Introduction

Organizations are increasingly seeking managerial strategies that will increase the effectiveness of information technology. The management of these information systems is a difficult task, challenged with balancing the divergent interests of many user groups in the face of rapid technological change. IS managers are confronted with the sometimes contradictory tasks of encouraging users to utilize newer technologies to derive additional benefits while ensuring that their use is cost-effective. In addition, such actions may diverge from an IS manager's personal agenda of increasing his or her span of control. It is, therefore, not surprising that IS departments are often unsuccessful in meeting the stated needs of their users.

### 4.2 Descriptive Results

The agency approach to EUC presented in this paper helps to explain the widespread occurrence of decentralized computing. In the absence of appropriate control mechanisms, end-users are likely to have opted for decentralized computing. Decentralization allows these users to make resource allocation decisions, including software and hardware acquisition decisions, and to develop implementation and operations schedules that are consistent with their self-interests. As discussed earlier, earlier IS environments were characterized by economies of scale and specialization in production that have decreased over time. An end-user manager would, therefore, have sought the decentralized solution at that point in time where the marginal costs of the externalities incurred plus the marginal costs that derive from loss of control over the information resource equal the decreasing marginal benefits of the economies of scale and specialization.

### 4.3 Prescriptive Results

Given the existence of decentralized computing and the trends in the technology, the theory provides insights into the appropriate division of IS activities between end-user departments and the IS department. It suggests that IS activities that experience large economies of scale or specialization relative to the cost of externalities should be centralized. These activities may include the use of large mainframes, telecommunications services and site licensing. On the other hand, if an activity is of relevance only to a single end-user department, the department manager should be free to determine how such a task is implemented. Perhaps more importantly, however, given the existence of interdependencies among most computing applications, a primary role of the IS organization must be to develop enforceable policies and standards that ensure that the costs to an end-user of developing computing applications or of using computing resources reflect the

true organizational costs, including the costs of externalities.

Agency theory highlights the possible differences in the goals of the IS manager, end-user managers and top-level managers, and it emphasizes the role of the differences in information possessed by each of these groups. These factors necessitate the implementation of control strategies that economize on agency costs. Agency theory suggests that these strategies focus on two major aspects of the control problem, the *informational aspects* and the *incentive aspects*.

One approach to addressing the existing information asymmetries is to increase the level of monitoring to improve the information that the principal possesses. However, the nature of information asymmetries in the IS context limits the value of monitoring as a means to reduce agency costs. For example, the output of an IS department is difficult to measure, the value of IS activities to users is similarly difficult to estimate, and the rapid pace of technological change makes it difficult to monitor the quality of decision-making by an IS manager.

An alternative approach is the use of incentive-compatible schemes that align the incentives of principals and agents. These schemes are designed in a manner such that agents are provided with incentives to provide accurate information. It is assumed that each agent possesses private information about his preferences and that he is self-interested. The objective is to achieve the optimal allocation of resources under these information asymmetries and goal incongruencies. These mechanisms typically involve a central planner who elicits information from agents and then determines a schedule of prices. Since it is virtually impossible to force agents to reveal their true valuations, the fee schedule is designed in a manner by which agents find it in their best interests to reveal their true valuations. A well-known example of such a scheme is the Clarke-Groves-Loeb (Clarke 1971; Groves and Loeb 1975) (hereafter CGL) tax mechanism. While there has been considerable focus on these schemes in the economics literature, relatively little work exists in the IS context. Work on the design of incentive-compatible schemes in the IS context is due primarily to the work of Mendelson and Whang.[8] Their work has focused primarily on the optimal allocation of mainframe resources under queuing delays. The CGL scheme could also be applied to other IS management issues. (See Appendix A for an example applying the CGL scheme to a software acquisition decision where there are multiple user departments and several competing software packages.)

Despite this literature in economics, the design of contracts to minimize agency costs that result from actions taken by the IS manager do not appear to have received any attention in the IS research and management literatures. There are a number of reasons why this might be the case. In addition to the lack of IS research attention to this area,

these schemes are sometimes difficult to implement. Moreover, IS activities are so varied that significant effort would be required to develop schemes that would address the multiple tasks.[9] Finally, the impacts of managerial actions in the IS context have not been well understood, and only now are managerial practices beginning to stabilize (Nolan 1979; Gurbaxani and Mendelson 1990). Indeed, there is still considerable variance in managerial opinion related to such issues as the choice of organization structure for the IS department (Swanson and Beath 1988) and even to the institution of chargeback systems (Allen 1987).

## 4.4 Summary

This paper has proposed that agency theory provides a useful framework within which to analyze managerial decision-making in the IS context. It has suggested that the widespread growth of EUC can be explained, in part, by the existence of agency costs in the IS environment. In addition, the agency model suggests that the use of incentive-compatible schemes can be used to decrease agency costs and improve the management of EUC. Development of formal hypotheses with which to empirically validate this approach would be a desirable next step. Future research using agency theory is likely to be successful both in explaining other observed phenomena and in developing better control mechanisms.

## 5. ACKNOWLEDGEMENTS

## 6. REFERENCES

Alavi, M. "End-User Computing: The MIS Managers' Perspective." *Information and Management*, Volume 8, 1985, pp. 171-178.

Alavi, M.; Nelson, R.; and Weiss, I. "Strategies for End-User Computing: An Integrative Framework." *Journal of Management Information Systems*, Winter 1987-88, Volume 4, Number 3.

Allen, B. "Make Information Services Pay Its Way." *Harvard Business Review*, Volume 65, Number 1, January-February 1987, pp. 57-65.

Arthur Andersen & Co. *The Changing Shape of MIS*. 1986.

Banker, R., and Kemerer, C. "Performance Evaluation of Information Systems Department Manager-Agents." MIT Sloan School Working Paper, November 1989.

Beath, C., and Straub, D. "Managing Information Resources at the Department Level: An Agency Perspective." *Proceedings of the Twenty-Second Hawaii International Conference on Systems Sciences*, Volume III, January 1989, pp. 151-159.

Benson, D. H. "A Field Study of End-User Computing: Findings and Issues." *MIS Quarterly*, Volume 7, Number 4, December 1983, pp. 35-45.

Clarke, E. H. "Multipart Pricing of Public Goods." *Public Choice*, Volume 11, 1971, pp. 17-33.

Cotterman, W., and Kumar, K. "User Cube: A Taxonomy of End Users." *Communications of the ACM*, Volume 32, Number 11, November 1989, pp. 1313-1320.

Davis, G., and Olson, M. *Management Information Systems: Conceptual Foundations, Structure, and Development.* New York: McGraw-Hill, 1985.

Dickson, G.; Leitheiser, R.; and Wetherbe, J. "Key Information Systems Issues for the 1980s." *MIS Quarterly*, Volume 8, Number 3, September 1984, pp. 135-162.

Eisenhardt, K. "Agency Theory: An Assessment and Review." *Academy of Management Review*, Volume 14, Number 1, January 1989, pp. 57-74.

Eisenhardt, K. "Control: Organizational and Economic Approaches," *Management Science*, Volume 31, Number 2, February 1985, pp. 134-149.

Fama, E., and Jensen, M. "Separation of Ownership and Control." *Journal of Law and Economics*, Volume 26, June 1983.

Gerrity, T. P., and Rockart, J. F. "End-User Computing: Are you a leader or a laggard?" *Sloan Management Review*, Volume 27, Number 4, Summer 1986, pp. 25-34.

Groves, T., and Loeb, M. "Incentives and Public Inputs." *Journal of Public Economics*, Volume 4, 1975, pp. 211-226.

Guimaraes, T. "The Benefits and Problems of User Computing." *Journal of Information Systems Management*, Fall 1984, pp. 3-9.

Guimaraes, T. "Personal Computing Trends and Problems: An Empirical Study." *MIS Quarterly*, Volume 10, Number 2, June 1986, pp. 179-187.

Gurbaxani, V., and Kemerer, C. "An Agent-Theoretic Perspective of the Management of Information Systems."

*Proceedings of the Twenty-Second Hawaii Conference on Systems Science*, Volume III, January 1989, pp. 141-150.

Gurbaxani, V., and Mendelson, H. "An Integrative Model of Information Systems Spending Growth." *Information Systems Research*, Volume 1, Number 1, March 1990, pp. 23-46.

Henderson, J. C. "Managing the IS Design Environment: A Research Framework." CISR Working Paper #158, Massachusetts Institute of Technology School of Management, September 1987.

Henderson, J. C., and Treacy, M. E. "Managing End-User Computing." *Sloan Management Review*, Winter, 1986.

Hirshleifer, J. *Price Theory and Applications.* Englewood Cliffs, New Jersey: Prentice Hall, 1980.

Huff, S. L.; Munro, M. C.; and Martin, B. H. "Growth Stages of End-User Computing." *Communications of the ACM*, Volume 31, Number 5, May 1988, pp. 542-550.

Jensen, M. "Organization Theory and Methodology." *The Accounting Review*, Volume LVIII, Number 2, April 1983.

Jensen, M. C., and Meckling, W. H. "Theory of the Firm: Managerial Behavior, Agency Costs and Ownership Structure." *Journal of Financial Economics*, Volume 3, October 1976, pp. 305-60.

Klepper, R. "An Agency Theory Perspective on Information Centers." *Proceedings of the Twenty-Third Annual Hawaii International Conference on System Sciences*, 1990, pp. 251-259.

Kriebel, C., and Moore, J. "Economics and Management Information Systems." In E. R. McLean (ed.), *Proceedings of the First International Conference on Information Systems*, Philadelphia, Pennsylvania, 1980, pp. 19-31.

Mendelson, H. *Economics of Information Systems Management.* Englewood Cliffs, New Jersey: Prentice-Hall, forthcoming 1990.

Mendelson, H. "Economies of Scale in Computing: Grosch's Law Revisited." *Communications of the ACM*, Volume 30, Number 12, December 1987, pp. 1066-1073.

Mendelson, H. "Pricing Computer Services - Queuing Effects." *Communications of the ACM*, Volume 28, Number 5, March 1985, pp. 312-321.

Mendelson, H., and Whang, S. "Optimal Incentive-Compatible Priority Pricing for the M/M/1 Queue." *Operations Research*, forthcoming, September/October 1990.

Munro, M. C.; Huff, S. L.; and Moore, G. "Expansion and Control of End-User Computing." *Journal of Management Information Systems*, Volume 4, Number 3, Winter 1987-88, pp. 6-27.

Nolan, R. L. *Management Accounting and Control of Data Processing.* New York: National Association of Accountants, 1979.

Ouchi, W. "A Conceptual Framework for the Design of Organizational Control Mechanisms." *Management Science*, Volume 25, September 1979, pp. 833-848.

Rivard, S., and Huff, S. "An Empirical Study of Users as Application Developers." *Information and Management*, Volume 8, 1985, pp. 89-102.

Rivard, S., and Huff, S. "Factors of Success for End-User Computing." *Communications of the ACM*, Volume 31, Number 5, May 1988, pp. 552-561.

Rivard, S., and Huff, S. "User Developed Applications: Evaluation of Success from the DP Department Perspective." *MIS Quarterly*, Volume 8, Number 1, March 1984, pp. 39-50.

Robey, D., and Zmud, R. "Research on End-User Computing: Theoretical Perspectives from Organization Theory." *Proceedings of IFIP Working Group 8.2 Conference on Desktop Information Technology*, Ithaca, New York, June 1989.

Rockart, J. F., and Flannery, L. S. "The Management of End-User Computing." *Communications of the ACM*, Volume 26, No.10, 1983, pp. 776-784; reprinted in J. Rockart and C. Bullen (eds.), *The Rise of Managerial Computing.* Homewood, Illinois: Dow Jones-Irwin, 1986, pp. 285-310.

Silberberg, E. *The Structure of Economics - A Mathematical Analysis.* New York: McGraw-Hill, 1978.

Swanson, E. B., and Beath, C. M. "Division of Labor and Departmentalization in Software Development and Maintenance." Working Paper, Anderson Graduate School of Management, University of California, Los Angeles, October 1988.

Whang, S. "Alternative Mechanisms of Allocating Computer Resources Under Queuing Delays." *Information Systems Research*, Volume 1, Number 1, March 1990, pp. 71-88.

Whang, S. "Pricing Computer Services: Incentive, Information and Queuing Effects." Unpublished Ph.D Dissertation, W.E. Simon Graduate School of Business Administration, University of Rochester, 1988.

Williamson, O. *The Economic Institutions of Capitalism.* New York: Free Press, 1985.

## 7. ENDNOTES

1. First most important in a list of twenty-two issues, Arthur Andersen (1986); second most important in a list of nineteen issues, Dickson, Leitheiser, and Wetherbe (1984).

2. See, for example, Benson (1983), Cotterman and Kumar (1989), Guimaraes (1986), Rivard and Huff (1984, 1985, 1988), Rockart and Flannery (1983).

3. Consistent with this approach, the behavioral assumption is made that all economic units act out of self-interest. While there are obvious divergences between the goals of individual actors within each of the three units and the goals of their managers, these are of secondary importance in the current analysis. Thus, for the purposes of exposition, the idealized assumptions are made that top management attempts to maximize the value of the firm, that end-users within a functional department attempt to maximize the "objective function" of the department head and that the central IS staff attempts to maximize the "objective function" of the IS manager. Henderson (1987) uses a similar approach in studying the IS design environment.

4. With the possible exception of organizations whose primary external product is information services, information systems are a support or staff function and, therefore, the IS department "works for," in an agency sense, the end-user departments, and not the other way around.

5. In particular, there were significant economies of scale in hardware technologies, as is embodied in Grosch's law (Mendelson 1987). Moreover, not only were the unit costs of computing high, but hardware capacity could only be purchased in large, discrete chunks.

6. For a more detailed discussion of these issues, see Gurbaxani and Kemerer (1989).

7. Note that information systems managers may also have their own preferences that may be driven by a desire to maximize their own market value by developing expertise on particular machine types.

8. See Mendelson (1985), Whang (1988), Mendelson and Whang (1990), and Whang (1990).

9. Banker and Kemerer (1989) have recently developed a model of multi-criteria performance evaluation in the IS context.

# APPENDIX A

## CLARKE-GROVES-LOEB TAX EXAMPLE

Consider the software acquisition decision where there are multiple user departments and several competing software packages. Each department manager is asked how much he or she is willing to pay for each package. For example, assume that there are three managers and three software packages, as shown in Table A.1.

### Table A.1 Differential Values of Software Packages

| Manager | Packages A | B | C | Tax |
|---------|-----------|-----|-----|-----|
| 1       | 50        | 20  | 0   | 30  |
| 2       | 0         | 60  | 20  | 0   |
| 3       | 40        | 0   | 50  | 30  |
| Total   | 90        | 80  | 70  | 60  |

The total value of each package is computed by summing over each managers' stated value for that package. The package that receives the highest total score is acquired. The key to ensuring that the managers reveal their true valuations is the chargeback mechanism. The difference between the values associated with any two packages is the dollar amount that a particular manager would be willing to pay to have the package with the higher value than the one with the lower value. By summing any particular column, total values for each package can be determined. In the example, package A is valued highest. The taxes can be computed by systematically determining the resulting outcome absent each one of the managers. These results are shown in Table A.2. For example, if manager 1 is excluded, package C would have been selected by a difference of $30 (70 − 40). Hence, manager 1 would be taxed $30, since it was due in part to his or her valuation that package A rather than package C was chosen. The surcharge, or tax, that manager 1 pays is the price for the privilege of determining the package eventually chosen. On the other hand, manager 2 would not be taxed, since package A would still be chosen without taking his or her preferences into consideration. Finally, manager 3 would be taxed $30 (80 − 50), since package B would be selected if manager 3 abstained from the process.

### Table A.2 Totals without Each Manager

| | Packages A | B | C |
|---|---|---|---|
| Without manager 1 | 40 | 60 | 70 |
| Without manager 2 | 90 | 20 | 50 |
| Without manager 3 | 50 | 80 | 20 |

# Empirical Tests of a Principal-Agent Model of the Investor-Investment Advisor Relationship

Joseph H. Golec*

## Abstract

This paper develops a specialized principal-agent model of the investor-investment advisor relationship and embeds the standard advisory compensation schedule in the model. Advisors are endowed with information-gathering abilities and investors are endowed with funds. Information-gathering services are traded indirectly through the investor's receipt of portfolio returns net of advisory fees. Model results show that the parameters of the compensation schedule are both a function of the idiosyncracies of an advisor's information services and the degree of risk sharing between the advisor and investor. Several predictions of the model are supported using data on mutual fund advisors. Unsupported predictions may be due to self-selection of advisors by risk tolerance.

## I. Introduction

The path-breaking work of Ross (1973) and Holmstrom (1979) has spawned numerous papers that characterize specific situations of economic exchange as principal-agent relationships. These include the shareholder-manager (Jensen and Meckling (1976)), mineral owner-extractor (Leland (1978)), government-contractor (Weitzman (1980)), issuer-investment banker (Baron (1982)), and investor-investment advisor (Starks (1987)) relationships, to name a few. While the principal-agent model is a popular way to describe many examples of economic interaction, very little empirical evidence exists to support it. Agrawal and Mandelker (1987), Tehranian, Travlos, and Waegelein (1987), Healy (1985), and Larcker (1983) link management decisions to management performance plans, but use indirect methods of analysis. This paper directly links the parameters of investment advisors' incentive contracts to the characteristics of their portfolio returns and empirically tests the relationships.

The primary result of the more interesting theoretical principal-agent models is an optimal compensation schedule that pays the agent a share of the

---

output. Recently, Holmstrom and Milgrom (1987) have shown that linear compensation schedules are optimal for realistic assumptions about the behavior of agents and the types of information observable to principals. Ramakrishnan and Thakor (1984) and Campbell and Kracaw (1985), (1987) use linear compensation schedules in their models but they do not fit a specific situation. Starks (1987) extends the approach by using a specific linear compensation schedule, namely, that used by investment advisors. She employs a general utility function that serves the purpose of her paper well, but does not allow the schedule parameters to be analyzed as readily as Ramakrishnan and Thakor and Campbell and Kracaw who use specific utility functions.

This paper blends both approaches. The linear fee schedule used by investment advisors is embedded in the structure of a simple principal-agent model and its optimal parameters are analyzed easily as functions of important economic variables. The specialized nature of the model yields comparative static results concerning the compensation schedule parameters. These detailed tests should be quite powerful.

This model differs from those of Starks, Ramakrishnan and Thakor, and Campbell and Kracaw in that the input of the agent is clearly specified as investment information rather than general effort. In addition, information (input) affects not only the expected value of portfolio return (output), but also affects return variance as a by-product; i.e., attempts to "beat the market" by taking positions based on information lead to portfolio-specific return variability. Campbell and Kracaw explicitly assume variance is unaffected and Starks and Ramakrishnan and Thakor assume this implicitly.

The paper is organized as follows. Section II describes the investor-investment advisor relationship and explains how it fits the assumptions of the principal-agent model. A model is constructed in Section III and some empirical implications are derived. Section IV describes the data, which include a sample of mutual fund investment advisors' fee schedules, and tests the implications of the model, focusing primarily on the cross-sectional variation of fee schedule parameters. Section V is a conclusion.

## II.    The Investor-Investment Advisor Relationship

The investor-investment advisor relationship can be characterized as a principal-agent relationship in which the investor (principal) hires an investment advisor (agent) to supply investment information (input) that affects the distribution of the investor's portfolio return (output). The investor does not receive the information directly but benefits through extra portfolio return. Because it is prohibitively costly for the investor to monitor the advisor, the advisor has no incentive to apply information to the portfolio unless the investor optimally constructs an incentive contract, subject to informational constraints, that provides the proper incentives.

For simplicity, assume that advisors are endowed with investment information or information-gathering abilities but have no capital to invest, and that investors are endowed with capital but no information. Investors believe that

information may be applied to investment portfolios such that the following return-generating process holds,

$$(1) \qquad \tilde{R}_p \;\; = \;\; \beta_p \tilde{M} + I + (I\delta)^{\frac{1}{2}} \tilde{\varepsilon}$$

where $\tilde{R}_p$ is the random gross portfolio return, $\tilde{M}$ is the random gross market return, $\beta_p$ (beta) is a real valued scalar, $I$ are the units of nonrandom return associated with the units of information applied to the portfolio, $\tilde{\varepsilon}$ is a unit of the random portfolio-specific return, and $\delta$ is the advisor's information ratio. Assume $E(\tilde{R}_p) = \bar{R}_p$, $E(\tilde{M}) = \bar{M}$, $E(\tilde{\varepsilon}) = 0$, $\text{Var}(\tilde{M}) = \sigma_M^2$, $\text{Var}(\tilde{\varepsilon}) = \sigma_\varepsilon^2$, and $\text{Cov}(\tilde{M}, \tilde{\varepsilon}) = 0$, where $E(\cdot)$, $\text{Var}(\cdot)$, and $\text{Cov}(\cdot, \cdot)$ are the expectation, variance, and covariance operators, respectively. $\tilde{M}$ and $\tilde{\varepsilon}$ are assumed to be normally distributed, hence $\tilde{R}_p$ is normally distributed.

The information ratio describes the tradeoff of additional units of portfolio-specific variance for each unit of nonrandom portfolio-specific return generated by the advisor's information. The investor must bear additional variability as a consequence of the advisor's actions, which concentrate investment in the securities the advisor believes offer superior returns.[1] By definition, a smaller information ratio implies that an advisor is better able to act on information while simultaneously insulating the portfolio from the effects of random portfolio-specific information.

If the investor decides to hire the advisor, then the mean of the return distribution shifts to the right by $I$ and its variance increases by $I\delta\sigma_\varepsilon^2$. Otherwise, the investor hires no advisor and holds a perfectly diversified portfolio that returns $\beta_p \bar{M}$. Obviously, investors who believe (1) holds do not believe markets are perfectly efficient with respect to the advisor's information. The market is not perfectly efficient in this model because of the assumption of asymmetric information, i.e., only advisors are endowed with information.

Investors can observe $\tilde{M}$, $\tilde{R}_p$, and $\beta_p$ without cost. But because $I$, $\delta$, and $\tilde{\varepsilon}$ are assumed to be prohibitively costly to observe, investors are unable to determine from a superior portfolio return whether an advisor has supplied much information or whether the random portfolio-specific return was unusually large and positive. Advisors have no incentive to disclose when luck is responsible, in fact, their incentives are quite the opposite. This asymmetry in information, the crux of the principal-agent problem, is especially relevant for mutual fund investors who typically invest relatively small amounts and thus may not find it cost-effective to monitor advisor performance.

## III.   The Model

Following the approach of Ramakrishnan and Thakor (1984), an incentive fee schedule is embedded into a basic one-period principal-agent model. A form

---

[1]This covers stock pickers and market timers. Stock pickers are usually more undiversified than market timers. Nevertheless, even highly diversified market timers add volatility to their portfolio returns by periodically shifting the beta of their portfolios above or below that of their target betas. Measured relative to the passive target portfolio, market timers will have additional portfolio-specific return variation.

of the following fee schedule is used by most open-end mutual fund advisors,[2]

$$(2) \qquad \phi(k_b, k_i, \tilde{R}_p, \tilde{R}_x, A) \quad = \quad k_b A \tilde{R}_p + k_i A(\tilde{R}_p - \tilde{R}_x),$$

where $A$ is the dollar amount of the investment, $\tilde{R}_p$ and $\tilde{R}_x$ are the gross returns on the managed portfolio and the index portfolio, respectively, and $k_b$ and $k_i$ are the base and incentive fee parameters, respectively. The index portfolio is assumed to be perfectly diversified so that $\tilde{R}_x = \beta_x \tilde{M}$, where $\beta_x$ is the beta of the index portfolio.

The terms " base fees" and "incentive fees" are somewhat misleading. Base fees also provide incentives to advisors to supply information and increase portfolio returns because they are paid at the end of each period, which means that the advisor earns a portion of both the initial assets invested and the return over the period. From a multiperiod perspective, base fees may provide risk-averse advisors with significant incentives to supply information, because as superior returns compound (assuming returns are not paid out to investors), assets and base fees grow. Superior returns may also attract assets from new investors or more assets from old investors. Only the one-period incentive is captured in the model.

The first component of the basic principal-agent model is the agent's utility function. Assume that, while their information endowments may not be identical, advisors have identical utility functions exhibiting risk aversion and their entire wealth is obtained from their fees, hence, they are undiversified. The representative advisor's certainty equivalent utility can be expressed as follows,

$$(3) \qquad U(I, \phi) \quad = \quad E(\phi) - \tau \sigma^2(\phi) \quad = \quad f(w, I),$$

where $E(\phi)$ and $\sigma^2(\phi)$ are the expected value and variance of the fee schedule, respectively, and $\tau$ is the positive risk-aversion parameter. $f(w, I)$ is a function that specifies the opportunity cost of the advisor's information, i.e., its value to the advisor in its next best alternative use. Many principal-agent models, including Ramakrishnan and Thakor's, assume that the effort of the agent is associated with disutility rather than an opportunity cost. In this model, the agent may be endowed with information, hence, its allocation entails no disutility. It is assumed for simplicity that $f_I = w$ and $w > 0$, where $f_I$ is the partial derivative of $f$ with respect to $I$. The specification of $f(w, I)$ may contain a constant term that represents a rent on an advisor's endowment of information.

The model's second component is the investor's objective function. Assume that investors have identical risk-averse preferences for wealth and that they are well diversified. The current value of an investment to a representative investor can be defined as[3]

$$V_0 \quad = \quad \alpha[E(\tilde{V}_1) - \overline{M} \text{Cov}(\tilde{V}_1, \tilde{V}_M)\{\sigma_M \sigma(\tilde{V}_M)\}^{-1}],$$

---

[2] A tally of the actively managed mutual funds appearing in *Moody's Bank and Finance Manual* (1985) shows that 476 of 548 funds (87 percent) used this type of schedule. Most investment advisors who manage pension fund accounts, individual accounts, corporate accounts, or other bond and stock investment accounts use similar asset-based fee structures. Admati and Pfleiderer (1988) have shown that fees based on the number of fund shares held by an investor allow the advisor to effectively charge for the response to the information he or she generates. Since all shares have the same value, this arrangement is comparable to asset-based fees.

[3] See Ramakrishnan and Thakor (1984).

where $\tilde{V}_1 = A\bar{R}_p - \phi$ is the terminal value of the investment, $\alpha$ is the risk-free discount factor, $\tilde{V}_M$ is the terminal value of the market portfolio, and $\sigma(\tilde{V}_M)$ is its standard deviation. The investor's objective is to maximize the current value of his or her investment through the choice of fee parameters, subject to the constraints that (3) holds, and that the advisor chooses the amount of information applied to the portfolio. The formal principal-agent problem is

$$\text{Maximize } V_0$$
$$k_b, k_i$$

(4)          $$\text{subject to } f(w, I) - U(I, \phi) \;=\; 0$$

(5)                    $$I \in \text{argmax } [f(w, I) - U(I, \phi)].$$

By using (1) and (2) and substituting (4) into $V_0$, the optimal sizes of the fee parameters are obtained from the following Lagrangian,

$$L \;=\; \alpha[A\bar{R}_p - \tau\text{Var}(\phi) - f(w, I) - \overline{M}\{(1 - k_b)A\beta_p - k_i A(\beta_p - \beta_x)\}]$$
$$+ \mu[w - k_b A - k_i A + A^2 \tau \delta \sigma_\varepsilon^2 (k_b^2 + 2k_b k_i + k_i^2)],$$

where $\text{Var}(\phi) = k_b^2 A^2 \beta_p^2 \sigma_M^2 + k_b^2 A^2 I \delta \sigma_\varepsilon^2 + k_i^2 A^2 (\beta_p - \beta_x)^2 \sigma_M^2 + k_i^2 A^2 I \delta \sigma_\varepsilon^2 + 2k_i k_b A^2 \beta_p (\beta_p - \beta_x)\sigma_M^2 + 2k_i k_b A^2 I \delta \sigma_\varepsilon^2$, and $\mu$ is the Lagrangian multiplier associated with (5). The first-order conditions for a maximum are listed in Appendix A. Note that the first and third terms of the fee variance represent fee variability due to the variance of market return transmitted through the base and incentive fees, respectively. Similarly, the second and fourth terms are due to portfolio-specific return variability. Finally, the last two terms represent the covariance between the two fees due to market- and portfolio-specific variance, respectively.

## A.    Model Solutions for the Fee Parameters

The structural solutions for $k_b$ and $k_i$ from the first-order conditions are

$$k_b \;=\; \frac{\mu - \alpha\overline{M}\beta_p - 2A\tau[\delta\sigma_\varepsilon^2(\mu - \alpha I) - \alpha\beta_p(\beta_p - \beta_x)\sigma_M^2]k_i}{2A\tau[\delta\sigma_\varepsilon^2(\mu - \alpha I) - \alpha\beta_p^2\sigma_M^2]}$$

and

$$k_i \;=\; \frac{\mu - \alpha\overline{M}(\beta_p - \beta_x) - 2A\tau[\delta\sigma_\varepsilon^2(\mu - \alpha I) - \alpha\beta_p(\beta_p - \beta_x)\sigma_M^2]k_b}{2A\tau[\delta\sigma_\varepsilon^2(\mu - \alpha I) - \alpha(\beta_p - \beta_x)^2\sigma_M^2]}.$$

Note that each parameter depends on the size of the other premultiplied by a term that represents the impact of the covariance between the fees.

One can see the interesting implications of the model more clearly by examining the reduced form solutions for $k_b$ and $k_i$ from the first-order conditions,

(6)          $$k_b \;=\; \frac{1}{2A\tau}\left[\frac{\overline{M}}{\beta_x \sigma_M^2} - \frac{(\beta_p - \beta_x)\mu}{\beta_x \delta\sigma_\varepsilon^2[\mu - \alpha I]}\right]$$

and

$$(7) \qquad k_i \;=\; \frac{1}{2A\tau} \left[ \frac{-\overline{M}}{\beta_x \sigma_M^2} + \frac{\beta_p \mu}{\beta_x \delta \sigma_\varepsilon^2 [\mu - \alpha I]} \right].$$

It is clear from (6) and (7) that the unrestricted solution implies that the parameters can take on negative or positive values. However, given that $\beta_p > 0$ and $\mu > 0$ (the shadow price of inducing more information from the advisor is positive[4]), then $[\mu - \alpha I]$ must also be positive for $k_i$ to be able to take on positive values.[5] As shown below, both parameters cannot be negative simultaneously and, indeed, the sum of the two parameters must be positive.

To obtain an intuitive feel for Equations (6) and (7), note that the fee schedule rewards the advisor for two services of value to the investor, systematic risk sharing and information supply. Indeed, the portfolio return, and thus the fee, is separable into market-based and information-based portions. Using (1) and (2), the expected fee can be written as

$$(8) \qquad E(\phi) \;=\; [\beta_p k_b + (\beta_p - \beta_x)k_i]A\overline{M} + [k_b + k_i]AI.$$

The first part of the fee rewards the advisor for bearing systematic risk: $k_b$ is applied to a return with a beta of $\beta_p$ and $k_i$ is applied to a return with a beta of $(\beta_p - \beta_x)$. The second part is information-based: each fee parameter is applied to a return with the same information-based component, hence, the parameters are equally weighted.

It is not surprising that the relative sizes of the parameters will depend upon their relative contributions in obtaining the efficient level of risk sharing and information supply. To see this, add (6) and (7) to obtain

$$(9) \qquad k_b + k_i \;=\; \frac{1}{2A\tau} \left[ \frac{\mu}{\delta \sigma_\varepsilon^2 [\mu - \alpha I]} \right].$$

Similarly, weight $k_b$ in (6) by $\beta_p$, and $k_i$ in (7) by $(\beta_p - \beta_x)$ and add to obtain

$$(10) \qquad \beta_p k_b + (\beta_p - \beta_x)k_i \;=\; \frac{1}{2A\tau} \left[ \frac{\overline{M}}{\sigma_M^2} \right].$$

Thus, ignoring the scalar $1/(2A\tau)$, the weighted sums depend upon different variables. $[k_b + k_i]$ depends on the marginal benefit-to-cost ratio that guides efficient information supply. An additional unit of information, $I$, produces benefits of $\mu$ for the investor, but costs $\delta \sigma_\varepsilon^2 [\mu - \alpha I]$. $[\beta_p k_b + (\beta_p - \beta_x)k_i]$ depends upon the marginal benefit-to-cost ratio of risk sharing. Sharing a unit of risk with the advisor nets the investor $\overline{M}$, but costs $\sigma_M^2$. The scalar, $1/(2A\tau)$, adjusts the risk costs so that they are applicable to an undiversified advisor with a risk-aversion coefficient of $\tau$ who manages a portfolio of $A$ dollars of assets.

Changes in $\beta_p$ only affect the relative attraction of $k_b$ and $k_i$ for risk sharing. Since the parameters are equally attractive with respect to information

---

[4]See Holmstrom (1979) for an explanation of why $\mu > 0$.
[5]Actually, $[\mu - \alpha I] > 0$ is nearly guaranteed by substituting for $k_b$ and $k_i$ in the third first-order condition $(L_\mu)$ in Appendix A and solving for $\mu$.

inducement, one would expect that their relative sizes will depend upon risk-sharing considerations. As $\beta_p$ increases, $k_i$ increases while $k_b$ decreases by equal amounts, since we know from (9) that $(k_b + k_i)$ is unaffected by a change in $\beta_p$.

The parameters in (6) and (7) are defined with respect to the reference point $\beta_x$, and its size relative to $\beta_p$. Suppose we assume $\beta_p = \beta_x$, then specialization takes place. $k_b$ is solely determined by the relative benefits and costs of systematic risk sharing. Nevertheless, it is applied to the gross portfolio return that includes both the market-based and the information-based portions, hence, it still affects information supply incentives. This explains the first term in large brackets in (7), which reduces $k_i$ by the amount of the effect that $k_b$ has on incentives. Besides this adjustment, $k_i$ is based exclusively on the relative costs and benefits of information supply, since if $\beta_p = \beta_x$, then $k_i$ is applied solely to the information-based portion of return and, thus, has no risk-sharing potential.

One might question why $k_i$ would be used at all when $\beta_p < \beta_x$, since, ignoring the information-based return, the advisor can expect to lose when his portfolio return is compared with that of the higher beta index portfolio. The answer is that $k_i$ can be used to attain an efficient amount of risk sharing when the $\beta_p$ that the investor prefers is larger than the beta that offers the efficient amount of risk sharing. Losses on the market-based portion of the incentive fee will be offset by perfectly correlated gains from the market-based portion of the base fee that would not have been forthcoming had the investor chosen a smaller $\beta_p$. Hence, $k_i$ can be used to reduce the market risk shared with the advisor. Conversely, when $\beta_p > \beta_x$, $k_b$ is not necessarily zero because it can be used to increase risk sharing.

## B.   Further Results when Institutional Restrictions Apply

Thus far, the signs of the fee parameters have not been restricted, although the model implies that their sum must be positive. Investment advisors, in general, and mutual fund advisors, in particular, employ positive fee parameters exclusively. The model guarantees that $k_b > 0$ if $\beta_p \leq \beta_x$, but if $\beta_p > \beta_x$ then the following condition must hold,

$$(11) \qquad \frac{\overline{M}}{(\beta_p - \beta_x)\sigma_M^2} \quad > \quad \frac{\mu}{\delta\sigma_\varepsilon^2[\mu - \alpha I]}.$$

To guarantee $k_i > 0$,

$$(12) \qquad \frac{\overline{M}}{\beta_p\sigma_M^2} \quad < \quad \frac{\mu}{\delta\sigma_\varepsilon^2[\mu - \alpha I]}.$$

Conditions (11) and (12) can be interpreted as comparisons of the relative marginal benefits and costs of risk sharing and information supply for portfolios with betas of $(\beta_p - \beta_x)$ and $\beta_p$, respectively. Condition (11) says that if the benefit-cost tradeoff for risk sharing when $k_i > 0$ (i.e., with beta risk of $[\beta_p - \beta_x]$) exceeds the benefit-cost tradeoff for information supply, then $k_b$ should be positive as well. Condition (12) says that if the benefit-cost tradeoff for risk

sharing when $k_b > 0$ is smaller than the tradeoff for information supply, then $k_i$ should also be positive.

It is not surprising that, holding $\beta_x$ fixed, as $\beta_p$ gets larger, it becomes more likely that $k_i$ will be positive and less likely that $k_b$ will be positive. This is because, as $\beta_p$ increases, the marginal costs of risk sharing increase faster than the marginal benefits, hence, risk sharing becomes less attractive for high beta portfolios. Since $k_i$ is applied to a portfolio return with the smaller beta of $(\beta_p - \beta_x)$, it becomes more attractive for risk sharing.

Although many pension fund advisors use both base and incentive fees, the majority of mutual fund advisors use only a base fee parameter (343 of 370 sample mutual funds). It is unlikely that the exact condition for $k_i = 0$ from (7) is met for such a large percentage of funds. In addition, most of the 343 funds claimed that their advisors make significant efforts to "beat the market," and, indeed, most chose portfolios with significant amounts of unsystematic risk. Therefore, it is important to consider the possibility that institutional constraints force $k_i = 0$ for many funds, with risk sharing and information inducement still being demanded by investors.[6]  Solving the model, in this instance, for $k_b$ yields

$$k_b \;=\; \frac{\mu - \alpha \overline{M} \beta_p}{2A\tau[\delta \sigma_\varepsilon^2(\mu - \alpha I) - \alpha \beta_p^2 \sigma_M^2]}.$$

Here, both risk sharing and information supply inducement are handled with one instrument. If there is no value to advisor information, then $\mu = 0$ and $I = 0$. In this case, aside from the scalar $1/(2A\tau)$, $k_b$ is determined solely by the ratio of risk-sharing benefits to costs. On the other hand, if $\beta_p = 0$, then $k_b$ is determined solely by the ratio of information benefits to costs. As long as $\beta_p > 0$ or $I > 0$, then $k_b > 0$. Since advisors must be paid positive amounts for both risk sharing and information supply, then when $\beta_p > 0$ and $I > 0$, other things equal, $k_b$ must be larger than if one of the terms is zero.

## C.    Empirical Implications of the Model

Because the following empirical analysis is cross-sectional, only $k_b$, $k_t$, $I$, $\delta$, $A$, and $\beta_p$ are discussed. The other explanatory variables—$\overline{M}$, $\sigma_M^2$, $\sigma_\varepsilon^2$, $\alpha$, and $\beta_x$—do not fluctuate cross-sectionally and, thus, will not affect the analysis. The model assumes that all advisors have identical $\tau$s. Unfortunately, risk preferences may fluctuate cross-sectionally and are not measurable. If risk aversion varies in a systematic way, then omitting $\tau$ from the empirical models may affect some of the results, the potential effects of which are discussed below.

The following hypotheses (comparative statics) are generated from the model's results. First, it is clear from (6) that $k_b$ is negatively related to $A$

---

[6]Up until the late 1960s, incentive fees were almost never used. Modigliani and Pogue (1975) note that by 1970, however, more than 100 mutual funds added incentive fees. But, by 1985, only 29 funds used incentive fees, as most returned to the traditional base fee. One reason for the drop in number of funds using incentive fees may be the 1970 amendments to the Investment Advisors Act of 1940 that placed restrictions on incentive fees. Another is that the act always required fees to be "reasonable," and incentive fees have been attacked as being unreasonable because bonuses could be earned even when investors suffered negative returns.

and $\beta_p$: the negative relationship between $k_b$ and $A$ captures the economies of scale in portfolio management, whereas the negative relationship between $k_b$ and $\beta_p$ demonstrates that as $\beta_p$ increases, $k_b$ becomes relatively less attractive for risk sharing. Because the sign of $(\beta_p - \beta_x)$ is indeterminate, the relationships between $k_b$ and $\delta$ and $I$ are also indeterminate.

Second, from (7), $k_i$ is positively related to $\beta_p$ and $I$, and negatively related to $\delta$ and $A$: $k_i$ and $\beta_p$ are positively related because, as $\beta_p$ increases, $k_i$ becomes relatively more attractive for risk sharing, whereas $k_i$ and $I$ are positively related because greater $I$ requires better compensation. The negative relationship between $k_i$ and $A$ is accounted for by economies of scale. $k_i$ is also negatively related to $\delta$; when $\delta$ increases, the risk cost associated with information supply increases, which decreases the optimal amount of $I$ and, hence, decreases compensation.

Finally, (9) offers some testable hypotheses. The sum of the parameters $(k_b + k_i)$ is unrelated to $\beta_p$, positively related to $I$, and negatively related to $\delta$ and $A$. Recall from (8) that the expected fee is composed of market- and information-based components and that $(k_b + k_i)$ is applied to the information-based return. Thus, $\beta_p$ has no impact on this portion of the return and should not affect $(k_b + k_i)$. More $I$ requires better compensation, hence $(k_b + k_i)$ and $I$ are positively related. A larger $\delta$ implies a larger risk cost per unit of $I$, so less is demanded, implying that $(k_b + k_i)$ should be smaller. $A$ and $(k_b + k_i)$ are negatively related as a consequence of economies of scale. Furthermore, holding $I$, $\delta$, and $A$ constant, $k_b$ and $k_i$ are negatively related. Indeed, $k_b$ should be larger when $k_i = 0$ than when $k_b > 0$ and $k_i > 0$.

## IV.    Data and Empirical Tests

The data base consists of 370 open-end mutual funds, 27 of which have fee schedules containing base and incentive fees.[7] The distribution of the stated fund objectives of the 370 (27) fund sample is 66 (7) aggressive growth, 151 (15) growth, 83 (3) growth and income, 49 (1) balanced, and 21 (1) special. For a fund to be included, it had to have monthly return data available over the six-year period from January 1, 1982, through December 31, 1987, invest at least partially in common stocks, and have information on its fee schedule available.[8]

The hypotheses concerning $A$ and $\beta_p$ are likely to be more reliably tested because these variables are easily measurable. Assets are measured in millions of dollars at the end of 1987. Beta is measured over 1982–1987 using the Capital Asset Pricing Model (CAPM) with the Standard and Poor's 500 stock index as the market proxy and the one-month Treasury bill return as the risk-free rate.

---

[7] Of 476 mutual funds listed in *Moody's Bank and Finance Manual* (1985), only 29 used incentive parameters during the sample period; two began operations in 1984 and, therefore, did not have enough return data to be included.

[8] Data were obtained from CDA Investment Technologies Inc.; *Moody's Bank and Finance Manuals,* published by Moody's Investor Services; *Investment Companies,* published by Wiesenberger Financial Services; and mutual fund prospectuses. Only 17 funds had to be eliminated because they had no fee schedule information.

*I* and $\delta$ are not easily measured; indeed, this is the essence of the principal-agent problem. Nevertheless, investors may use imperfect measures to help determine fee parameters. The CAPM facilitates the use of Jensen's (1968) alpha as a proxy for *I*. Alpha could be negative even though the advisor applies information to the portfolio if random portfolio-specific returns are negative. $\delta$ is proxied by the standard deviation of a fund's portfolio-specific returns; that is, the standard deviation of the residuals from the CAPM (henceforth, called "residual"). This is not a pure measure of $\delta$ because it combines $\delta$ and *I* multiplicatively. Unfortunately, the effects of *I* cannot be eliminated by dividing residual by alpha because nonsensical numbers result when alpha is negative.

Fee parameters are measured as annual percentage rates and have been checked for changes. Few funds changed their fee schedules; however, for those that did, the parameters were averaged. None of the funds using an incentive parameter changed fee schedules, although two have eliminated their incentive parameters since 1987. Some funds explicitly state in their schedules how $k_b$ decreases for larger amounts of assets. For example, a fund might pay its advisor 0.5 percent of assets for the first \$100 million and 0.4 percent on assets above \$100 million. If the fund actually has \$200 million of assets, then $k_b$ is calculated as 0.45 percent. Although typical for very large funds, the majority of funds in the sample do not make $k_b$ conditional on assets. Finally, for one fund, $k_b = 0$ because it is an index fund and hired no advisor.

Panels I and II of Table 1 give an idea of the magnitudes of the fee parameters. The sample of 370 mutual funds is split into one sample of 343 funds that use only the base fee (Panel I) and another with 27 funds that use both base and incentive fees (Panel II). $k_b$ ranges from 0 to 1.7 percent with an average of 0.60 percent. $k_t$ ranges from 0.5 to 3.3 percent with an average of 1.57 percent for the 27 funds that include an incentive fee parameter. The incentive fee is paid based on performance relative to an index, typically the S&P 500 stock index.[9]

The sample of funds using incentive fees has a slightly larger mean alpha, although both samples of funds exhibited negative mean alphas during the sample period. The incentive funds managed portfolios with more assets, larger betas, and larger standard deviations of residual returns.

Most tests are performed using ordinary least squares regression analysis. The general form for the cross-sectional regressions is

$$\text{Fee Parameter}_j = b_0 + b_1 \text{Log}(\text{Assets}_j) + b_2 \text{Beta}_j + b_3 \text{Residual}_j + b_4 \text{Alpha}_j + e_j,$$

where $b_0$ is the intercept; $b_1$, $b_2$, $b_3$, and $b_4$ are regression coefficients; $e_j$ is a random error term with $E(e_j) = 0$, $\text{Var}(e_j) = \sigma_e^2$, and $\text{Cov}(e_i, e_j) = 0, i \neq j$; and $j$ denotes the $j$th mutual fund. The natural logarithm of assets is used because of the wide range of assets for the funds (see Table 1) and because the relationships between assets and fee parameters may not be linear. If advisors receive

---

[9]Although many mutual funds never intend to hold most or even some of the S&P 500 stocks, they use the S&P 500 index because it is well known to investors and its return performance is representative of that of publicly traded stocks. See for example, the December 16, 1985, prospectus of the Explorer II Fund.

TABLE 1

Descriptive Statistics for the Variables Used in Regression Tests of the Model

| | $k_b$ | $k_i$ | Assets | Beta | Residual | Alpha |
|---|---|---|---|---|---|---|
| | | | Regression Variables | | | |
| *Panel I. Sample of 27 mutual funds that have fee schedules with $k_b > 0$ and $k_i > 0$* | | | | | | |
| Mean | 0.56 | 1.57 | 1184.00 | 1.00 | 2 59 | −2.14 |
| Std. Dev. | 0.22 | 0.86 | 2384.70 | 0.22 | 1.90 | 4.52 |
| Minimum | 0.15 | 0.50 | 19.00 | 0.54 | 0.98 | −12.30 |
| Maximum | 1.00 | 3.30 | 11590.00 | 1.61 | 10.91 | 5.30 |
| *Panel II. Sample of 343 mutual funds that have fee schedules with $k_b > 0$ and $k_i = 0$* | | | | | | |
| Mean | 0.60 | — | 475.75 | 0.89 | 2.42 | −2.20 |
| Std. Dev. | 0.22 | — | 769.65 | 0.26 | 1.70 | 5.90 |
| Minimum | 0.00 | — | 0.70 | 0.30 | 0.00 | −38.10 |
| Maximum | 1.70 | — | 5284.00 | 2.21 | 15.58 | 12.10 |
| *Panel III. Correlation matrices[a] for the variables by sample* | | | | | | |
| $k_b$ | 1.00 | — | −0.38 | 0.13 | 0.19 | −0.30 |
| $k_i$ | 0.74 | 1.00 | — | — | — | — |
| Log (Assets)[b] | −0.44 | −0.50 | 1.00 | −0.09 | −0.17 | 0.46 |
| Beta | 0.50 | 0.37 | −0.18 | 1.00 | 0.21 | −0.65 |
| Residual | 0.32 | 0.37 | −0.19 | 0.15 | 1.00 | −0.37 |
| Alpha | −0.30 | −0.16 | 0.48 | −0.70 | −0.32 | 1.00 |

Means, standard deviations, minimums, and maximums of the variables used in the regression tests of the model. The tests are performed using data on two samples of mutual funds. The funds' base and incentive fee parameters ($k_b$ and $k_i$) are measured in annual percentage rates. Fund assets are measured in millions of dollars at the end of 1987. Their alphas, betas, and residuals are calculated over 1982–1987 with monthly data using the Capital Asset Pricing Model (CAPM), where the Standard and Poor's 500 stock index is the market proxy and the one-month Treasury bill return is the risk-free rate. Alphas are compounded to annual percentage rates. Each fund's residual is the standard deviation of the residuals from the CAPM.

[a] The correlations for the sample of 27 funds, where $k_b > 0$ and $k_i > 0$ appear below the diagonal, and those for the sample of 343 funds, where $k_b > 0$ and $k_i = 0$ appear above the diagonal.

[b] This is the natural logarithm of assets that is used in the regressions instead of assets.

incentives through increased assets, as discussed in Section II, fee parameters will not be reduced commensurately to hold the dollar level of fees constant. Instead, as assets increase, the parameters may be decreased by progressively smaller amounts.

The relatively high correlation of alpha with the rest of the independent variables in Panel III of Table 1 poses a problem for accurate estimation of the coefficients in the regressions. The remedy used here is to regress alpha on the other independent variables and use the residuals from this regression as the measure of information, rather than alpha itself. The corrected alpha has no correlation with the other independent variables.

Presentation of the regression results follows the order of the hypotheses in the previous section. Panel I of Table 2 includes regressions for $k_b$, $k_i$, and $k_b + k_i$, since the sample of 27 funds has $k_b > 0$ and $k_i > 0$; the 343-fund sample has $k_b > 0$ and $k_i = 0$, hence, Panel II contains a single regression for $k_b$. All of

the regression residuals were tested for heteroskedasticity of unspecified form using the Breusch-Pagan statistic. In this case, it has a chi-squared distribution with degrees of freedom equal to four. The largest observation for the statistic was 3.93, which is statistically insignificant.

TABLE 2

The Relationships between the Fee Parameters or Parameter Combinations and the Explanatory Variables: Cross-Sectional OLS Regression Results for the Model,

Fee Parameter$_j$   =   $b_0 + b_1$Log(Assets$_j$) $+ b_2$Beta$_j + b_3$Residual$_j + b_4$Alpha$_j + e_j$,

(where fee parameter$_j$ is fund $j$'s base fee parameter ($k_b$), incentive fee parameter ($k_i$), or the two combined ($k_b + k_i$), measured in annual percentage rates)

| Dependent Variable | Coefficients[a] | | | | | $R^2$ | F-Stat | $N$[b] |
|---|---|---|---|---|---|---|---|---|
| | $b_0$ | $b_1$ | $b_2$ | $b_3$ | $b_4$ | | | |
| Panel I. Sample of 27 mutual funds that have fee schedules with $k_b > 0$ and $k_i > 0$ | | | | | | | | |
| $k_b$ | 0.362 | −0.046* | 0.416* | 0.022 | 0.028* | 0.52 | 5.98* | 27 |
| | (1.62) | (2.18) | (2.68) | (1.26) | (2.29) | | | |
| $k_i$ | 1.548* | −0.217* | 1.007# | 0.114# | 0.157* | 0.63 | 9.49* | 27 |
| | (2.06) | (3.04) | (1.92) | (1.90) | (3.83) | | | |
| $k_i + k_b$ | 1.910* | −0.263* | 1.423* | 0.137# | 0.185* | 0.66 | 10.54* | 27 |
| | (2.18) | (3.16) | (2.33) | (1.96) | (3.86) | | | |
| Panel II. Sample of 343 mutual funds that have fee schedules with $k_b > 0$ and $k_i = 0$ | | | | | | | | |
| $k_b$ | 0.769* | −0.051* | 0.064 | 0.015* | −0.004 | 0.17 | 18.00 | 343 |
| | (13.40) | (7.14) | (1.48) | (2.26) | (1.38) | | | |

Fund assets are measured in millions of dollars at the end of 1987 and transformed to natural logarithms. Fund alphas, betas, and residuals are calculated over 1982–1987 with monthly data using the Capital Asset Pricing Model (CAPM), where the Standard and Poor's 500 stock index is the market proxy and the one-month Treasury bill return is the risk-free rate. Alphas are compounded to annual percentage rates. Each fund's residual is the standard deviation of the residuals from the CAPM.

[a] T-statistics in parentheses.

[b] $N$ is the number of observations.

* Significant at the 5-percent level

# Significant at the 10-percent level.

Judging from the $F$-statistics, all of the regressions are statistically significant. The model has identified variables that account for a significant proportion of variation in the parameters. The $R^2$s for the sample of 27 funds are quite high for cross-sectional regressions. The modest $R^2$ for the larger sample is still reasonable for a cross-sectional regression.

Several predictions were supported in Table 2. The first regression in Panel I shows that $k_b$ is negatively related to $A$, but is not negatively related to $\beta_p$. In the second regression, $k_i$ is positively related to $\beta_p$ and $I$ and negatively related to $A$, but is not negatively related to $\delta$. In the third regression, the sum of the parameters $k_b + k_i$ is negatively related to $A$ and positively related to $I$, but is not negatively related to $\delta$ or unrelated to $\beta_p$.

To decide whether $k_i$ and $k_b$ are negatively related, we need to control for the common effects that the independent variables $A$, $\beta_p$, $\delta$, and $I$ have on both

222

of them. Therefore, we use the residuals from the first two regressions in Panel I for $k_b$ and $k_i$, respectively, and regress the $k_b$ residuals on the $k_i$ residuals to obtain

$$k_b = -0.00 + 0.13^*k_i \quad R^2 = 0.20 \quad F\text{-statistic} = 6.41^*.$$
$$(t = 0.00)(t = 2.53)$$

Therefore, it appears that after eliminating the common effects, $k_b$, and $k_i$ are not negatively related, contrary to the prediction.[10]

Finally, a comparison of the average $k_b$ between samples is informative. A simple difference in means test using the mean $k_b$s in Panels I and II of Table 1 shows that their difference is statistically insignificant. On the other hand, this test does not control for differences between the samples in assets, alpha, beta, and residual. A more accurate test compares the intercept of the $k_b$ regression for the sample where $k_i > 0$ (Panel II of Table 2), to that of the $k_b$ regression for the sample where $k_i = 0$. After controlling for variation in $A$, $\beta_p$, $\delta$, and $I$, $k_b$ is larger on average when $k_i = 0$, as expected.

Several of the model's predictions are supported. In particular, the economies of scale in managing large amounts of assets is captured in the negative relationships between assets and the parameters. Confirmation of this relationship for $k_b$ is not surprising, given that some funds explicitly include decreases in $k_b$ for larger assets. The larger negative coefficient on assets in the $k_i$ regression is more surprising, since none of the funds explicitly provided for decreases in $k_i$ as assets increased.

Information (alpha) is positively related to the parameters when $k_i > 0$ and $k_b > 0$; more information supplied implies more compensation. Indeed, judging by the coefficients on alpha in the first two regressions of Panel I, alpha's effect on $k_i$ is greater than on $k_b$, as would be expected if $k_i$ is used more heavily than $k_b$ to reward information supply.

Predictions of negative relationships between the fee parameters and the risk variables ($\beta_p$ and $\delta$) or between the fee parameters themselves are not supported. Instead, in most cases, significant positive relationships are observed. Cross-sectional variation in the unobservable advisor risk-aversion parameter, $\tau$, may be the prime cause (although variation in $\mu$ may also be important). Because $\tau$ enters the denominators of all of the expressions for $k_b$ and $k_i$, as $\tau$ gets smaller (larger), both fee parameters should get larger (smaller). If advisors with relatively small $\tau$ choose to manage high-risk portfolios, and those with large $\tau$ choose low-risk portfolios, $\tau$ will be correlated with $\beta_p$ and $\delta$. Therefore, $\beta_p$ and $\delta$ may be picking up the positive effects that $\tau$ has on the parameters. In addition, the excluded effects of $\tau$ end up in the residuals of the $k_b$ and $k_i$ regressions in Panel I of Table 2. Hence, the residuals are positively rather than negatively correlated.

Of course, alpha and residual may not be good proxies; in particular, residuals may be used by investors as a proxy for information rather than the information ratio. Work by Trueman (1988) suggests that investment managers

---

[10]This positive relationship exists if the unadjusted $k_b$ and $k_i$ are used in the regression or if $A$, $\beta_p$, $\delta$, and $I$ are also included in the regression with the unadjusted parameters.

may trade excessively in order to give investors the appearance that they are trading on information. If investors use residual return variability as a measure of the trading activity of a manager and, hence, of his or her information, the hypotheses offered for alpha should also apply to residual. Alpha is expected to be positively related to the parameters and, as Table 2 shows, the coefficient on residual μ is positive in each regression.

## V.    Conclusions

This paper shows that the principal-agent model, sometimes thought to be too abstract for applied research, can be used to construct detailed analyses of specific economic relationships. Empirical test results for a specialized principal-agent model of the investor-advisor relationship are mixed. On the one hand, variables specified by the model consistently enter the regressions with significant coefficients and, more often than not, with the predicted sign. Moreover, for the sample of funds that used base and incentive fees, the regressions explained a large proportion of the variation in the fee parameters. On the other hand, some variables have coefficients with signs opposite from those predicted, which may be due to unmeasurable cross-sectional variation in advisor risk aversion.

The one-period model ignores the incentives that risk-averse advisors have to provide input in order to retain investors over many periods. This may be a moot point since the model shows that incentive fees may be predicated on risk sharing rather than pure incentive considerations. In any case, this study shows how the principal-agent model can be used to gain insight into the parameters of compensation schedules used for investment advisors, and can be adapted easily to other types of fee schedules.

## Appendix A

The Lagrangian shown in the text can be differentiated with respect to the choice variables to obtain,

$$(A\text{-}1) \quad L_b = 2A^2\tau[\delta\sigma_\varepsilon^2(\mu - \alpha I) - \alpha\beta_p^2\sigma_M^2]k_b - \mu A + \alpha\overline{M}A\beta_p$$
$$+ 2A^2\tau[\delta\sigma_\varepsilon^2(\mu - \alpha I) - \alpha\beta_p(\beta_p - \beta_x)\sigma_M^2]k_i = 0,$$

$$(A\text{-}2) \quad L_i = 2A^2\tau[\delta\sigma_\varepsilon^2(\mu - \alpha I) - \alpha(\beta_p - \beta_x)^2\sigma_M^2]k_i - \mu A + \alpha\overline{M}A(\beta_p - \beta_x)$$
$$+ 2A^2\tau[\delta\sigma_\varepsilon^2(\mu - \alpha I) - \alpha\beta_p(\beta_p - \beta_x)\sigma_M^2]k_b = 0,$$

$$(A\text{-}3) \quad L_\mu = w - k_b A - k_i A + A^2\tau\delta\sigma_\varepsilon^2[k_b^2 + 2k_b k_i + k_i^2] = 0,$$

where $L_b$, $L_i$, and $L_\mu$ are the partial derivatives of the Lagrangian with respect to $k_b$, $k_i$, and $\mu$.

Differentiating the first-order conditions again yields

$$L_{bb} = 2A^2\tau[\delta\sigma_\varepsilon^2(\mu - \alpha I) - \alpha\beta_p^2\sigma_M^2],$$

$$L_{ii} = 2A^2\tau[\delta\sigma_\varepsilon^2(\mu - \alpha I) - \alpha(\beta_p - \beta_x)^2\sigma_M^2],$$

$$L_{bi} = L_{ib} = 2A^2\tau[\delta\sigma_\varepsilon^2(\mu - \alpha I) - \alpha\beta_p(\beta_p - \beta_x)\sigma_M^2], \text{ and}$$

$$L_{\mu b} = L_{\mu i} = -A + A^2\tau\delta\sigma_\varepsilon^2[2k_b + 2k_i].$$

The sufficient second-order condition for a maximum with one constraint is that the determinant of the following bordered hessian be positive.

$$H = \begin{vmatrix} L_{bb} & L_{bi} & L_{\mu b} \\ L_{ib} & L_{ii} & L_{\mu i} \\ L_{\mu b} & L_{\mu i} & 0 \end{vmatrix} > 0.$$

This condition reduces to $\alpha\beta_x^2 > 0$, hence, a maximum is obtained.

# References

Admati, A. R., and P. Pfleiderer. "Selling and Trading on Information in Financial Markets." *American Economic Review,* 78 (May 1988), 96–103.

Agrawal, A., and G. N. Mandelker. "Managerial Incentives and Corporate Investment and Financing Decisions." *Journal of Finance,* 42 (Sept. 1987), 823–837.

Baron, D. P. "A Model of the Demand for Investment Banking Advising and Distribution Services for New Issues." *Journal of Finance,* 37 (Sept. 1982), 955–976.

Campbell, T. S., and W. A. Kracaw. "The Market for Managerial Labor Services and Capital Market Equilibrium." *Journal of Financial and Quantitative Analysis,* 20 (Sept. 1985), 277–297.

_____. "Optimal Managerial Incentive Contracts and the Value of Corporate Insurance." *Journal of Financial and Quantitative Analysis,* 22 (Sept. 1987), 315–328.

Healy, P. "The Effects of Bonus Schemes on Accounting Decisions." *Journal of Accounting and Economics,* 7 (April 1985), 85–107.

Holmstrom, B. "Moral Hazard and Observability." *Bell Journal of Economics,* 10 (Spring 1979), 74–91.

Holmstrom, B., and P. Milgrom. "Aggregation and Linearity in the Provision of Intertemporal Incentives." *Econometrica,* 55 (March 1987), 303–328.

Jensen, M. C. "The Performance of Mutual Funds in the Period 1945–1964." *Journal of Finance,* 23 (May 1968), 386–416.

Jensen, M. C., and W. H. Meckling. "Theory of the Firm: Managerial Behavior, Agency Costs, and Ownership Structure." *Journal of Financial Economics,* 3 (Oct. 1976), 305–360.

Larcker, D. "The Association between Performance Plan Adoption and Corporate Capital Investment." *Journal of Accounting and Economics,* 5 (April 1985), 209–232.

Leland, H. E. "Optimal Risk Sharing and the Leasing of Natural Resources, with Application to Oil and Gas Leasing on the OCS." *The Quarterly Journal of Economics,* 92 (Aug. 1978), 413–437.

Modigliani, F., and G. Pogue. "Alternative Investment Performance Fee Arrangements and Implications for SEC Regulatory Policy." *Bell Journal of Economics,* 6 (Spring 1975), 127–160.

*Moody's Bank and Finance Manual.* New York: Moody's Investor Services (1985).

Ramakrishnan, R. T. S., and A. V. Thakor. "The Valuation of Assets under Moral Hazard." *Journal of Finance,* 39 (March 1984), 229–238.

Ross, S. A. "The Economic Theory of Agency: The Principal's Problem." *American Economic Review,* 63 (May 1973), 134–139.

Starks, L. "Performance Incentive Fees: An Agency Theoretic Approach. *Journal of Financial and Quantitative Analysis,* 22 (March 1987), 17–32.

Tehranian, H.; N. G. Travlos; and J. F. Waegelein. "The Effects of Long-Term Performance Plans on Corporate Sell-Off-Induced Abnormal Returns." *Journal of Finance,* 42 (Sept. 1987), 933–942.

Trueman, B. "A Theory of Noise Trading in Securities Markets." *Journal of Finance,* 43 (March 1988), 83–95.

Weitzman, M. L. "Efficient Incentive Contracts." *Quarterly Journal of Economics,* 94 (June 1980), 719–730.

# Agency Theory and Its Application to Small Firms: Evidence from the Swedish Venture Capital Market

## Hans Landström

The research in small firms financing is characterized by a lack of a theoretical framework. One basic assumption in this study is that agency theory can provide an essential framework to explain the interaction between informal and formal venture capitalists and their portfolio firms. Five hypotheses generated from agency theory are formulated and tested on 62 firms backed by informal venture capitalists and 145 firms backed by formal venture capitalists. The theoretical conclusion is that agency theory does not provide a satisfactory framework to explain either the informal venture capitalist's, nor the formal venture capitalist's relationship to their portfolio firms. Therefore, more exploratory research must be done to develop a theory of finance which will be applicable in the small firms situation.

## I. INTRODUCTION

The research in small firms financing is characterized by a lack of a theoretical framework. The theory of modern corporate finance is developed mainly with large firms in mind. Ang [1] emphasizes that small firms have unique characteristics. For example, the owners have undiversified personal portfolios, the first generation owners are entrepreneural and prone to risk taking, the management team is not complete, the small firms experience high cost of market and institutional imperfections, and relationships with stockholders are less formal. These differences between large and small firms could generate a different set of financial problems in small firms, or cause the entrepreneurs in the small firms to look at the same financial problems in a different manner.

Against this background, it is important to develop or find theories which are valid in the small firms situation. One of the assumptions in this study is that agency theory can provide a helpful framework to explain the cooperation between venture capitalists and their portfolio firms.

Hans Landström ● Halmstad University, Box 823, S-301 18, Halmstad, Sweden.

Why is agency theory interesting in this respect? The following reasons can be mentioned:

- Agency theory has been developed on a relatively high level of abstraction, but the theory is to a limited extent tested in empirical situations.
- Researchers in the field of agency theory regard small firms as one area where studies can provide most of the leverage for agency theory [6].
- Agency costs are thought to be a major impediment to small firms in their attempts to obtain external financing [2, 13, 21].
- Some of the assumptions in the agency theory sound intuitively appealing and relevant in the small firms situation. For example, the high level of asymmetric information [2], and the fact that entrepreneurs are not motivated to disclose information due to fears that it might be used against them, so-called information impactedness [14].
- Studies of the venture capital market in the US have shown that the agency theory could be a helpful theory to explain the interaction between formal venture capitalists and entrepreneurs [22], and that monitoring financial and operation performance is one of the most time-consuming activities of formal venture capitalists [15].

The aim of the study is to test the applicability of agency theory to small firms. The empirical material in the study is derived from the Swedish venture capital market.

The term "venture capital" is often misused, and there is no universal definition [16, 24]. However, essential to the definition is that the venture capitalists provide risk capital (equity and near-equity capital), the investments are made in small unlisted firms, and the commitments are for a limited period of time.

The concept "informal venture capitalists" will be defined as external private individuals who provide risk capital directly to small unlisted firms. The definition of informal venture capitalists in this study is broader than the definition of informal investors used in studies in the US [9, 25] and the UK [17]. Furthermore, the concept "formal venture capitalists" refers to companies (with no strong connection to private individuals' capital) which provide risk capital to small unlisted firms.

The paper is organized into seven sections. In section II a review of the agency theory is presented. Section III presents the hypotheses generated from agency theory. Section IV describes the data collecting process and the variables used in the study. Some characteristics of the firms surveyed are

presented in section V. In section VI the empirical findings are presented and in the final section some theoretical conclusions are drawn.

## II. AGENCY THEORY

In this section the general assumptions in the agency theory will be presented, and the external investors' possibilities to monitor the entrepreneur's behavior will be described.

In general, agency theory is related to the problem that occurs when cooperating parties have different goals and a division of labor. Specifically, the agency theory focuses on the relationship in which one or more persons (the principal(s)) engage another person (the agent) to perform some work on their behalf [13]. The basic premise of agency theory is that both principals and agents are assumed to be rational economic-maximizing individuals. Therefore, the separation of ownership and control will result in decisions by the agent which are not always in the principal's best interest and there will arise costs (agency costs) of bringing the agent's behavior into line. For example, costs arise which are incurred by the principals when monitoring and controlling the behavior of the agent (so-called monitoring costs), and costs incurred by the agent in demonstrating compliance with the wishes of the principal (so-called bonding costs).

The unit of analysis in the agency theory is the contract between the principal and agent. These contracts (written and unwritten) specify the rights of the agent, performance criteria on which agents are evaluated, and the payoff functions they face [8]. Especially, there are two problems that the agency theory tries to solve. The first is the problem that arises when the goals of the principal and agent conflict and it is difficult or expensive for the principal to verify what the agent is actually doing. The second is the problem that arises when the principal and agent have different attitudes toward risk, which can lead to different preferred actions.

The agency theory has its roots in information economics, and the theory has developed along two lines; positivist and principal-agent research [12]. The two approaches share a common unit of analysis and use the same agency cost minimizing tautology, but differ in their mathematical strictness. Positivist research is less mathematical and more empirically oriented than principal-agent research. The positivist researchers have focused mainly on the principal-agent relationship between owners and managers of large corporations [e.g., 7, 8, 13], whereas principal-agent researchers are concerned with a general theory of the principal-agent relationship [e.g., 4, 11].

One of the core issues in the agency theory concerns the principals' possibilities to monitor the agent's behavior. Monitoring refers to the

principals' ability to determine whether the agents have lived up to the provisions of the contract and to prevent the agent's misuse of assets due to conflicts of interest. In Jensen-Meckling's [13] definition, monitoring refers to more than just measuring or observing the behavior of the agent. It also includes efforts to "control" the behavior of the agent through budget restrictions, operating rules, etc.

In the case were the principal does not have complete information about the agent's behavior, as in the case of external investors, two options exist [5]; to put the agent's behavior under surveillance (e.g., through reporting procedures, and board of directors), or to reward the agent based on outcomes (e.g., profitability).

Following such reasoning Ouchi [19, 20] suggests two underlying monitoring strategies. The strategy can be either behavior or outcome based. The behavior-based strategy refers to an agreement between the principal and the agent which concerns a certain behavior that in some way will be rewarded, whereas outcome-based strategy refers to the principal's measurement of certain outcomes and the reward will be based on this measurement. According to Ouchi [19], the choice between the strategies depends on two dimensions; knowledge of transformation process and availability of output measures. To use a behavior-based strategy, that is, to continuously observe the agent's behavior, the principal requires a causal knowledge of what is required to attain a desired outcome. When the principal uses outcome-based strategy, for example, to measure the agent's attained results, the transformation process need not be known at all, but a reliable and valid measure of the desired outputs must be available.

## III. HYPOTHESES

In this section the hypotheses generated from agency theory will be presented. The hypotheses emanate from Ouchi's reasoning regarding different monitoring strategies. As a monitoring variable I will use the concept "active involvement" which includes a high frequency of contacts between the investor and entrepreneur, and more operational involvement by the investor in the portfolio firm.

The characteristics of the portfolio firm can be expected to affect the venture capitalist's way of monitoring the entrepreneur. According to Ouchi [19] a more behavior-based strategy will be used in situations where the availability of output measures are low. This is the situation in highly innovative firms and in young firms. Furthermore, high innovation firms may involve technology that is not well understood by the venture capitalist. The information asymmetries increase the threat of opportunism which leads

to more active involvement. In young firms the risk of failure is high. This also leads to a need for frequent contacts and more operational involvement. In accordance with this reasoning the following hypotheses can be formulated:

H1: The higher the innovation level of the firm, the more the venture capitalist will rely on more active involvement.

H2: The younger the firm, the more the venture capitalist will rely on more active involvement.

Environmental characteristics will also affect the monitoring behavior. As the environment becomes more variable, information will have to be processed more frequently. Under conditions of high variability, performance evaluation becomes more difficult [13] and more evaluation mechanisms are likely to evolve [3]. Therefore, it is reasonable to expect that variability in the environment will force a more active involvement. The hypothesis is:

H3: The more variable the environment becomes, the more the venture capitalist will rely on more active involvement.

The effects of ownership on managerial incentives are one of the core issues of agency theory. One suggestion is that an entrepreneur who owns a large share of the firm will require little monitoring, because his incentives will be in line with those of outside owners. Furthermore, as the ownership of the outside investors increases, the need for monitoring will increase. The reasons for this are; (I) the risk that the entrepreneur will consume the firm's resources will increase, and (II) the investors exposure to business risks increases as the equity stake increases. As a result, the venture capitalist will take a more active role in the firm when his ownership level is high. This reasoning leads to the following hypothesis:

H4: The lower the relative ownership level of the entrepreneur, the more the venture capitalist will rely on more active involvement.

The venture capitalist's knowledge about the portfolio firm's transformation process, here defined as the knowledge of the firm's market and technology, will affect the monitoring behavior. According to Ouchi [19] a more behavior-based strategy will be used in situations where the principal's

**Table 1**
**Hypotheses in the Study**

|  |  | Monitoring Variables | |
|---|---|---|---|
|  |  | Frequency of Contacts | Operational Involvement |
| H1 | Innovation Level | + | + |
| H2 | Age of the Firm | + | + |
| H3 | Environment Variability | + | + |
| H4 | Entrepreneur's Ownership | − | − |
| H5 | Venture Capitalist's Knowledge | + | + |

knowledge of the tranformation process is good. This will give rise to the following hypothesis:

H5: The more the venture capitalist knows about the portfolio firm's transformation process, the more the venture capitalist will rely on more active involvement.

In Table 1 the hypothesized impact of the context on the venture capital—entrepreneur relationships is summarized.

## IV. METHOD

In this section the data collecting process and the variables used in the study will be presented. Furthermore, the limitations of the study will be discussed. The study is based on two surveys, one of firms backed by informal venture capitalists and one survey of firms backed by formal venture capitalists.

### Survey of firms backed by informal venture capitalists

This part of the study was carried out during the spring of 1991, and is based on a survey of manufacturing and technology-based firms in Sweden. Three geographic areas in southern Sweden and 11 science parks were selected for the study. The criteria for the sample in the three geographic areas were manufacturing firms with up to 100 employees. The sample frame was composed of a random sampling from the data base of the Postal Office (PAR). The science parks were studied through a full-scale survey of the firms located at the science parks.

In total 1,258 firms were included in the sample frame. The questionnaire was mailed to the CEO's of the firms, with a reminder via telephone after two weeks and a postal reminder after additionally one week. Of the 1,258 firms, 47 claimed that they were not independent juridical firms,

not manufacturing firms or firms with over 100 employees. Sixteen questionnaires were sent back by the postal services, and it was assumed that those firms had gone out of business. Thirty-one firms reported that they had discontinued their operation or had gone into bankruptcy. The effective sample frame was thus 1,164 firms, and of these 627 firms were not heard from, 32 firms sent back incomplete questionnaires or questionnaires that were not filled in, and 505 sent back questionnaires that were completely filled in. Thus, the response rate is 505/1,164, or 43%.

The results show *inter alia* that banks, as might be expected, are the most commonly used external investor, followed by supplier and leasing/factoring companies. It is interesting to note that informal venture capitalists are used in 62 or 12% of the firms, which is suprisingly high. The analyses in this paper are based on those 62 firms which have informal venture capitalists as (part) owner.

## Survey of firms backed by formal venture capitalists

In this survey the data is based on a questionnaire sent to CEO's in firms backed by formal venture capitalists in Sweden. The survey was carried out in the spring of 1990. The sample frame was designed from the venture capital data base of the Swedish National Board for Industrial and Technical Development, and from annual reports of venture capital companies. In total 536 portfolio firms were traced. Commitments which involved one venture capital company investing in another venture capital company were disregarded. The fact that one portfolio firm may have several venture capital companies as (part) owner has also been taken into account. To offset the risk that the questionnaire might be filled in by representatives of the venture capital company, portfolio firms with the same address as the venture capital company were excluded. The final list of portfolio firms used in the survey included 380 firms.

The questionnaire was mailed to the portfolio firms, with a reminder after three weeks. Of the 380 firms, answers were obtained from 183. Of these, there were 17 firms with more than 100 employees. Thus, the effective sample frame was 363 firms. 21 questionnaires were returned "blank" and 145 questionnaires could be used for analyses. The percentage of answers was thus 145/363, or 40%.

## Variables used in the surveys

The variables that have been used in the study were operationalized in the following way (Table 2).

### Table 2
### Overview of Variables Used in the Study

| Variables | Operationalization |
| --- | --- |
| Innovation Level | Biomodel scale (0 = Old product on the market and 1 = New product on the market) |
| Age of the Firm | Year of Start |
| Environment Variability | Five point scale (1 = Small changes to 5 = Large changes) in the dimensions, market, technology, competition and supliers |
| Entrepreneur's Ownership Share | Percent |
| Venture Capitalist's Knowledge | Five point scale (1 = Very limited extent to 5 = Very large extent) in the dimensions, market and technology |
| Frequency of Contacts | Five point scale (1 = Almost never, 2 = When needed, 3 = Monthly, 4 = Weekly and 5 = Daily) |
| Operational Involvement | Seven point scale (1 = No active cooperation, 2 = Economic reports, 3 = Work on board, 4 = Ad hoc when needed, 5 = Continuous informal contacts, 6 = Involvement in operation [part time] and 7 = Involvement in operation [full time]) |

## Limitations of the study

There are several factors which potentially restrict the conclusions which may be drawn. First, one such limitation is the size of the samples, including the survey of firms backed by informal venture capitalists as well as the survey of firms backed by formal venture capitalists. Larger samples would have been preferred for statistical analyses and generalization purposes. Secondly, the most serious limitation refers to the operationalizing of the variables. Some of the variables are measured through single item measures due to the desire to get an acceptable response rate. This can be discussed since the contents in some of the variables are more comprehensive than what can be included in a single item measure. Furthermore, the construct validity for separate variables can be discussed. For example, the assumption behind the variable "venture capitalist's knowledge about the portfolio firm's transformation process" is that this knowledge will be reflected in the venture capitalist's provision of resources in the dimensions, market, and technology. Finally, the broad definition of "informal venture capitalists" makes it difficult to compare the results in this study with studies of informal investors in the US and the UK.

## V.  SOME CHARACTERISTICS OF THE FIRMS SURVEYED

In this section the firms backed by formal and informal venture capitalists'
will be described.

### The character of the firms

The survey of firms backed by formal venture capitalists showed that
44% of the firms were started during the 1980's. The average number of
employees were 32 (median 25 employees), and 22% of the firms had less than
10 employees. Twenty-two percent of the firms had a principal product that
was "completely new on the market" when it was launched.

The average number of owners in the firms were 2.7 owners. The average
share owned by the CEO was 16% (median 0%), and the formal venture
capitalist was majority owner in 59% of the cases.

Corresponding results in the survey of firms backed by informal venture
capitalists showed that 69% of the firms were started during the 1980's. The
firms in this survey are smaller. The average number of employees was 16
employees (median 11 employees), and 44% had less than 10 employees.
Among the informal venture capitalists' portfolio firms 49% stated that the
firm had a principal product that was "completely new on the market" when
it was launched.

On average there were 9.5 owners in the firms. The average share owned
by the CEO was 30% (median 26%), and the informal venture capitalists were
majority owners in 49% of the firms.

Due to differences in the sample frame it is difficult to make any
conclusions regarding the investment patterns between informal and formal
venture capitalists.

### The Cooperation Between the Venture
### Capitalists and the Portfolio Firms

Of course, the differences in the sample frame also influence the
possibilities to make comparisons between the formal and informal venture
capitalists' way of cooperating with their portfolio firms. However, in both
cases the relationship is formed between parties that are rather close. The
distance between the formal venture capitalist and the portfolio firm was
less than 50 km in 51% of the cases. Corresponding results for the informal
venture capitalists were 60% of the cases.

The frequency of contacts between the venture capitalists and the
portfolio firms are rather high for both formal and informal venture
capitalists. Forty-eight percent of the respondents in the formal venture

**Table 3**
**Provision of Resources**

| | Average Values on a Five Point Scale | |
| --- | --- | --- |
| | Formal Venture Capitalists | Informal Venture Capitalists |
| Acting as a Sounding Board | 3.1 | 3.7 |
| Wider Range of Contacts | 2.9 | 3.4 |
| Facilitated Contacts with Interested Third Parties | 2.7 | 3.1 |
| Professionalizing of the Portfolio Firm | 2.8 | 3.1 |
| Financial Expertise | 3.4 | 3.1 |
| Expertise in Negotiating and Contract-Making | 2.6 | 3.1 |
| Market Expertise | 1.8 | 2.6 |
| Technological/Production Expertise | 1.4 | 2.5 |

capitalists survey and 53% in the informal venture capitalists survey indicated contacts some time/s every day or week.

There are some differences in the way of organizing the cooperation between the venture capitalists and the portfolio firms. The formal venture capitalists seem to rely more on financial reports and consultancy work in the portfolio firms, whereas informal venture capitalists are more actively involved in operations. Both formal and informal venture capitalists work actively on board meetings and by informal contacts with the entrepreneurs.

Apart from capital, the venture capitalists provide different kinds of expertise. This is primarily in the form of acting as a sounding board, professionalizing of the portfolio firm, financial expertise, and expertise in negotiating and contract-making. A comparison between formal and informal venture capitalists shows that the informal venture capitalists provide expertise to a larger extent on almost every studied variable (see Table 3).

The portfolio firms' expectations of the cooperation with the venture capitalists have in many cases been fulfilled. Seventy-four percent of the CEO's in firms backed by informal venture capitalists are of the opinion that their expectations had been realized to a large or very large extent. Corresponding results for the entrepreneurs in firms backed by formal venture capitalists were 54%.

## VI.   EMPIRICAL RESULTS

In this section the five hypotheses are tested. The effects of the independent variables on the venture capitalists frequency of contacts and operational involvement are analyzed.

## Table 4
### The Impact of Independent Variables on Frequency of Contacts

*Dependent Variable: Frequency of Contacts*

| Independent variables | Hypothesis | Prediction | Formal Venture Capitalists | | Informal Venture Capitalists | |
|---|---|---|---|---|---|---|
| | | | Beta Value | Significance | Beta Value | Significance |
| Innovation level | H1 | + | −0.19 | | −0.54 | Sign. |
| Age of the firm | H2 | + | 0.00 | | 0.01 | |
| Environment variability | H3 | + | 0.02 | | −0.35 | Sign. |
| Entrepreneur's ownership | H4 | − | −0.01 | | 0.01 | |
| Venture capitalist's knowledge | H5 | + | 0.23 | *** | 0.29 | ** |
| $R^2$ (adj) | | | | 0.19 | | 0.33 |
| F | | | | 5.04*** | | 5.56*** |
| n | | | | 89 | | 46 |

*Notes:* Level of Significance
+ $p < 0.10$
\* $p < 0.05$
\*\* $p < 0.01$
\*\*\* $p + 0.001$

By way of introduction the correlations between the variables used as independent variables should be examined. The primary interest is to examine the extent to which multicollinearity can be expected to confound the results of the regression analyses conducted to test the hypotheses.

The result of the correlations shows that none of correlations are above 0.50. The highest value in both surveys is between the variables innovation level and age of the firm ($r = 0.23$ in the formal venture capitalist survey, and $r = 0.32$ in the informal venture capitalist survey). This indicates that the variables are tapping different aspects of the venture capitalist—entrepreneur relationship, and it appears that multicollinearity should not be a serious threat to the regression analysis.

A basic assumption in the study is that venture capitalists attempt to manage the agency risks inherent in a particular firm through the level of their involvement. Table 4 presents the results of regressing the five independent variables against frequency of contacts in the venture capitalist—entrepreneur dyad.

The results in Table 4 show that only one of the five hypotheses is supported. As predicted, the frequency of contacts increases when the venture

**Table 5**
**The Impact of Independent Variables on Operational Involvement**

| Independent Variables | Hypothesis | Prediction | Formal Venture Capitalists | | Informal Venture Capitalists | |
|---|---|---|---|---|---|---|
| | | | Beta Value | Significance | Beta Value | Significance |
| Innovation level | H1 | + | 0.30 | | −0.60 | Sign. |
| Age of the firm | H2 | + | −0.01 | | 0.01 | |
| Environment variability | H3 | + | −0.19 | | −0.11 | |
| Entrepreneur's ownership | H4 | − | 0.00 | | −0.00 | |
| Venture capitalist's knowledge | H5 | + | 0.52 | *** | 0.39 | * |
| $R^2$ (adj) | | | | 0.11 | | 0.13 |
| F | | | | 3.33** | | 2.46* |
| n | | | | 93 | | 51 |

*Notes:* Level of Significance
　　　+ $p < 0.10$
　　　* $p < 0.05$
　　　** $p < 0.01$
　　　*** $p + 0.001$

capitalist has more knowledge about the portfolio firm's transformation process (H5). None of the other hypotheses are supported. On the contrary, for the informal venture capitalists there is a significant relationship between higher innovation level, respectively higher environment variability, and lower frequency of contacts (H1 and H3). Also, the results do not give support for the hypotheses concerning a positive relationship between young firms and a high frequency of contacts (H2), or a negative relationship between the entrepreneur's ownership level and the frequency of contacts (H4).

The predicted direction of the hypotheses is identical when operational involvement is used as a dependent variable. Table 5 presents the results of the regression analysis testing the hypotheses regarding the effects of the independent variables on the venture capitalist's operational involvement in the portfolio firms.

The results for the hypotheses regarding the venture capitalists operational involvement in the portfolio firms show similar results as for the frequency of contacts. Only the variable "venture capitalist's knowledge" (H5) is significantly related to operational involvement. The variables, age of the firm (H2), environment variability (H3), and entrepreneur's ownership

(H4) seem to contribute little to explaining variations in the operational involvement by the venture capitalists in their portfolio firms. It is interesting to note that high innovation level is related to more operational involvement by the formal venture capitalists, but less operational involvement by the informal venture capitalists (H1).

To summarize, it appears that only one variable generated from agency theory helps to explain the active involvement by the venture capitalists in their portfolio firms. More knowledge from the venture capitalist about the portfolio firm's transformation process seems to support higher frequency of contacts and operational involvement. It is also interesting to note the difference between informal and formal venture capitalists in the treatment of portfolio firms with a high innovation level. In these situations, the informal venture capitalists seem to be considerably more passive than the formal venture capitalists.

The weak relationships between the variables generated from agency theory and the venture capitalists' involvement in their portfolio firms may be explained by the differences in the venture capitalists' ownership level. Therefore, a comparison was made between those firms with majority ownership by venture capitalists against those with minority ownership. However, the result shows no major differences between minority and majority owned portfolio firms. Thus, the results do not strongly support the assumption that the venture capitalists take a more active role in the portfolio firms when their ownership level is high.

## VII.   CONCLUSIONS

This section summarizes the conclusions emerging from the regression analyses testing of the formulated hypotheses. The section is divided into two subsections; theoretical implications and discussion.

### Theoretical implications

This study used agency theory as the theoretical framework to study the relationship between venture capitalists and entrepreneurs. It is interesting to note that the agency theory does not seem to provide a satisfactory explanation for the venture capitalist's interaction with the entrepreneurs. This holds true for formal as well as informal venture capitalists. The results are contrary to what was expected, and the results are especially interesting for the formal venture capitalists, since the assumptions in agency theory could be expected to be more valid in the relationship between formal venture capitalists and their portfolio firms. It is possible that the agency theory is

not valid in the relationship between venture capitalists and their portfolio firms due to the following reasons:

I.   The agency theory is based on the assumption that both principals and agents are rational economic-maximizing individuals. This does not hold for the entrepreneur or the informal venture capitalist. Studies have shown that entrepreneurs are often driven by other than purely economical motives. Also, studies of informal venture capitalists in the US and the UK show that they do not always see the monetary rewards as the most essential.

II.  The agency theory assumes that the principal building control mechanism is to prevent opportunistic behavior from the agent, which implies a "negative" relationship between the principal and agent. The relationship between the venture capitalist and entrepreneur usually has a more "positive" character, where the interaction is based on support and mutual trust. In many cases the control mechanism functions as a dysfunctional factor with lowering trust between the venture capitalist and entrepreneur, which impedes open communication, etc.

III. The agency theory assumes that there is an information asymmetry between the principal and agent which facilitates the agent's opportunistic behavior. The negotiations between the venture capitalist and entrepreneur, and the personal relationship between them can result in less information asymmetries and less opportunistic behavior, and therefore substitute monitoring solutions.

My conclusion is that the agency theory is not applicable in the interaction between venture capitalists and entrepreneurs. More exploratory research must be done to develop a theory of finance which will be applicable in the small firms situation.

## Discussion

Finally, some reflections regarding the differences between informal and formal venture capitalists. As mentioned earlier, the differences in the sample frame imply that a complete comparison cannot be obtained. However, some observations can be made.

The results in the study indicate that informal venture capitalists have a tendency to invest in young firms and technology-based firms to a larger extent than formal venture capitalists. This corresponds with results in the US [10, 23, 25] and in the UK [18].

In addition to the differences in the investment pattern there seems to exist differences in the ownership structure in the portfolio firms. Firms backed by informal venture capitalists seem to have a weaker ownership structure, with a large number of owners and where each informal venture capitalist has a small ownership share in the firm. This implies that the cooperation must be based on a mutual trust between the venture capitalist and entrepreneur. On the other hand, firms backed by formal venture capitalists seem to have a stronger ownership structure with few owners and where the venture capitalist in many cases is the majority owner. The consequence is that the formal venture capitalist to a larger extent can rely on its ownership power in their cooperation with the portfolio firm.

Of course, the differences in investment pattern and ownership structure influence the conditions of the cooperation between the venture capitalists and portfolio firm. The results in the study indicate that the informal venture capitalists provide more expertise to their portfolio firms compared to the formal venture capitalists. However, the impression is that informal and formal venture capitalists react differently to changes in the portfolio firms. The formal venture capitalists react quicker and more resolute. They are more objective in their judgements, and economic motives guide their decisions. Informal venture capitalists react less powerfully, they are more subjective in their judgements, and they must to a larger extent rely on the entrepreneur's statements and actions. Therefore, the informal venture capitalists appear to provide less assistance with short term changes and/or problems in the portfolio firms than the formal venture capitalists.

## Acknowledgments

## REFERENCES

[1]  Ang, J.S. 1991. Small Business Uniqueness and the Theory of Financial Management. *Journal of Small Business Finance* 1(1): 1-13.

[2]  Barnea, A., R. Haugen, and L. Senbet. 1981. Market Imperfections, Agency Problems and Capital Structure: A Review. *Financial Management* (Summer): 7-22.

[3]  Bourgeois, L.J. 1985. Strategic Goals, Perceived Uncertainty, and Economic Performance in Volatile Environments. *Academy of Management Journal* 28: 548-573.

[4]  Demski, J., and G. Feltham. 1978. Economic Incentives in Budgetary Control Systems. *Accounting Review* 53:336-359.

[5]   Eisenhardt, K.M. 1985. Control: Organizational and Economic Approaches. *Management Science* 31(2):134-149.

[6]   Eisenhardt, K.M. 1989. Agency Theory: An Assessment and Review. *Academy of Management Review* 14(1):57-74.

[7]   Fama, E. 1980. Agency Problems and the Theory of the Firm. *Journal of Political Economy* 88(2):288-307.

[8]   Fama, E., M. and Jensen. 1983. Separation of Ownership and Control. *Journal of Law and Economics* 26:301-325.

[9]   Gaston, R.J. 1989. *Finding Private Venture Capital For Your Firm: A Complete Guide.* Wiley: New York.

[10]  Haar, N.E., J. Starr, and I.C. MacMillan. 1988. Informal Risk Capital Investors: Investment Patterns on the East Coast of the USA. *Journal of Business Venturing* 3:11-29.

[11]  Harris, M., and A. Raviv. 1979. Optimal Incentive Contracts with Imperfect Information. *Journal of Economic Theory* 20:231-259.

[12]  Jensen, M.C. 1983. Organization Theory and Methodology. *Accounting Review* 56(2):319-338.

[13]  Jensen, M.C., and W.H. Meckling. 1976. Theory of the Firm: Managerial Behavior, Agency Costs and Ownership Structure. *Journal of Financial Economics* 3:305-360.

[14]  Kaplan, R.S., and A.A. Atkinson. 1989. *Advanced Management Accounting.* Englewood Cliffs, NJ: Prentice-Hall.

[15]  MacMillan, I.C., D.M. Kulow, and R. Khoylian. 1988. Venture Capitalists' Involvement in their Investments: Extent and Effect In *Frontiers of Entrepreneurship Research,* edited by B.M. Kirchhoff, W.A. Long, W.E. McMullan, K.H. Vesper, and W.E. Wetzel, Jr., 303-323. Wellesley, MA: Babson College.

[16]  Maier, J., and D. Walker. 1987. Role of Venture Capital in Financing Small Business. *Journal of Business Venturing* 2:207-214.

[17]  Mason, C., and R. Harrison. 1990. Informal Risk Capital: A Review and Research Agenda. Venture Finance Research Project, Working Paper No. 1, University of Southampton/University of Ulster at Jordanstown.

[18]  Mason, C., R. Harrison, and J. Chaloner. 1991. Informal Risk Capital in the UK: A Study of Investor Characteristics, Investment Preferences and Investment Decision-making. Venture Finance Research Project, Working Paper No. 2, University of Southampton/University of Ulster at Jordanstown.

[19]  Ouchi, W. 1977. The Relationship Between Organizational Structure and Organizational Control. *Administrative Science Quarterly* 22(March):95-113.

[20]  Ouchi, W. 1979. A Conceptual Framework for the Design of Organizational Control Mechanisms. *Management Science* 25:833-848.

[21]  Pettit, R.R., and R.F. Singer. 1985. Small Business Finance: A Research Agenda. *Financial Management* 14(3):47-60.

[22]  Sapienza, H.J. 1989. *Variations in Venture Capitalist—Entrepreneur Relation: Antecedents and Consequences.* University of Maryland at College Park.

[23]  Tymes, E.R., and O.J. Krasner. 1983. Informal Risk Capital in California. In *Frontiers of Entrepreneurship Research,* edited by J.A. Hornaday, J.A. Timmons, and K.H. Vesper, 347-368. Wellesley, MA: Babson College.

[24]  Wan, V. 1991. Australian Venture Capital Market Revisited. *Technovation* 11(6):327-337.

[25]  Wetzel, W.E. 1981. Informal Risk Capital in New England. In *Frontiers of Entrepreneurship Research,* edited by K.H. Vesper, 217-245. Wellesley, MA: Babson College.

# STAKEHOLDER–AGENCY THEORY

CHARLES W. L. HILL
THOMAS M. JONES

*School of Business Administration, University of Washington*

## ABSTRACT

Taking agency theory and stakeholder theory as points of departure, this article proposes a paradigm that helps explain the following: (1) certain aspects of a firm's strategic behaviour; (2) the structure of management–stakeholder contracts; (3) the form taken by the institutional structures that monitor and enforce contracts between managers and other stakeholders; and (4) the evolutionary process that shapes both management–stakeholder contracts and the institutional structures that police those contracts.

## INTRODUCTION

Over the last decade, agency theory has emerged as the dominant paradigm in the financial economics literature (Jensen and Meckling, 1976; Ross, 1973). As developed in that literature, agency theory has been primarily concerned with the relationship between managers and stockholders. However, recently authors in the management field have begun to explore the implications that agency theory might have for the disciplines of organizational behaviour, organizational theory, and strategic management (*e.g.*, Eisenhardt, 1985, 1988, 1989; Kosnik, 1987). One area that remains relatively unexplored concerns the ability of agency theory to explain the nature of the implicit and explicit contractual relationships that exist between a firm's stakeholders. In addition to managers and stockholders, stakeholders include employees, customers, suppliers, creditors, communities, and the general public. The agency theory view of the firm as a nexus of contracts between resource holders (stakeholders) suggests that this may be a promising avenue for investigation.

Taking agency theory and stakeholder theory as points of departure, the purpose of this article is to propose a paradigm that helps explain the following: (1) certain aspects of a firm's strategic behaviour; (2) the structure of management–stakeholder contracts; (3) the form taken by the institutional structures that monitor and enforce contracts between managers and other stakeholders; and (4) the evolutionary process that shapes both management–stakeholder contracts and the institutional structures that

police those contracts. Like agency theory, this paradigm suggests that the firm can be seen as a nexus of contracts between resource holders. Unlike agency theory, the paradigm encompasses the implicit and explicit contractual relationships between *all* stakeholders. Stated simply, the resultant model is a generalized theory of agency: one of *stakeholder–agency*.

Although similar to agency theory in many respects, stakeholder–agency theory is based on assumptions concerning market processes that are substantially different from those underlying the finance version of agency theory. The result is a paradigm whose predictions are not always consistent with those of agency theory. While agency theory operates on the assumption that markets are efficient and adjust quickly to new circumstances, here the existence of short to medium-run market inefficiencies are admitted. The result is the introduction of power differentials into the stakeholder–agent equation. Although the idea of power differentials is at variance with the traditional agency approach, in our view the approach developed here increases the explanatory power of the paradigm.

## AGENCY THEORY

An agency relationship is defined as one in which one or more persons (the principal(s)) engages another person (the agent) to perform some service on their behalf which involves delegating some decision-making authority to the agent (Jensen and Meckling, 1976; Ross, 1973). The cornerstone of agency theory is the assumption that the interests of principles and agents diverge. According to agency theory, the principal can limit divergence from his/her interests by establishing appropriate incentives for the agent, and by incurring monitoring costs designed to limit opportunistic action by the agent. Further, it may pay the agent to spend resources (bonding costs) to guarantee that he/she will not take certain actions that would harm the principal, or to ensure that the principal will be appropriately compensated if he/she does take such action. That is, the agent may incur *ex-ante bonding costs* in order to win the right to manage the resources of the principal. Despite these devices, it is recognized that some divergence between the agent's actions and the principal's interests may remain. Insofar as this divergence reduces the principals's welfare, it can be viewed as a *residual loss*.

The sum of the principal's monitoring expenditures, the agent's bonding expenditures, and any remaining residual loss are defined as agency costs. Further, agency theory asserts that natural selection processes favour governance structures that economize on agency costs (Fama and Jensen, 1983; Jensen, 1983). By governance structures, agency theorists mean the mechanisms that police the explicit and implicit contracts between principals and agents (Demsetz, 1983; Fama, 1980; Fama and Jensen, 1983). These include the structure of law governing corporate behaviour and its attendant legal apparatus, monitoring mechanisms (such as the board of directors), and enforcement mechanisms (such as the market for corporate control and the managerial labour market).

Although applied primarily to the stockholder–manager relationship, Jen-

sen and Meckling (1976) argue that agency theory 'will lead to a rich theory of organizations which is now lacking in economics and the social sciences generally' (p. 309). Jensen and Meckling view the implicit contract between stockholders and managers as just one of the nexus of contracts that form the legal fiction known as the modern corporation. Other contracts that could be considered within an agency framework include those between managers and the various primary interest groups of the firm or stakeholders.

## STAKEHOLDERS

The term stakeholders refers to groups of constituents who have a legitimate claim on the firm (Freeman, 1984; Pearce, 1982). This legitimacy is established through the existence of an exchange relationship. Stakeholders include stockholders, creditors, managers, employees, customers, suppliers, local communities, and the general public. Following March and Simon (1958), each of these groups can be seen as supplying the firm with critical resources (contributions) and in exchange each expects its interests to be satisfied (by inducements). Stockholders provide the firm with capital. In exchange, they expect the firm to maximize the risk-adjusted return on their investment. Creditors provide the firm with finance and in exchange expect their loans to be repaid on schedule. Managers and employees provide the firm with time, skills, and human capital commitments. In exchange, they expect fair income and adequate working conditions. Customers supply the firm with revenues and expect value for money in exchange. Suppliers provide the firm with inputs and seek fair prices and dependable buyers in exchange. Local communities provide the firm with locations, a local infrastructure, and perhaps favourable tax treatment. In exchange, they expect corporate citizens who enhance and/or do not damage the quality of life. The general public, as tax payers, provides the firm with a national infrastructure. In exchange, they expect corporate citizens who enhance and/or do not damage the quality of life and do not violate the rules of the game established by the public through their legislative agents.

### Specific Asset Investments

Stakeholders differ with respect to the size of their stake in the firm. The magnitude of an individual actor's stake is a function of the extent to which that actor's exchange relationship with the firm is supported by investments in specific assets (Williamson, 1984). Following Williamson (1984, 1985), by specific assets we mean assets that cannot be redeployed to alternative use without a loss of value. For example, employees with general-purpose skills and knowledge can leave the firm and be replaced without productive loss to either the worker or the firm (assuming efficient labour markets). In such cases, their 'stake' in the firm is low. Alternatively, employees with skills that are uniquely tailored (specialized) to the requirements of the firm cannot leave without bearing substantial exit costs in the form of the lower rent stream that their skills can earn in the next best application. The 'stake' of such employees in the firm is high. This distinction is important: compared to

actors with a low stake in the firm, actors with a high stake will demand more comprehensive incentive mechanisms and governance structures in order to safeguard their asset-specific investments in the firm.

## The Unique Role of Management

Whatever the magnitude of their stake, each stakeholder is a part of the nexus of implicit and explicit contracts that constitutes the firm. However, as a group, managers are unique in this respect because of their position at the centre of the nexus of contracts. Managers are the only group of stakeholders who enter into a contractual relationship with all other stakeholders. Managers are also the only group of stakeholders with *direct* control over the decision-making apparatus of the firm (although some stakeholders, and particularly the suppliers of capital, have indirect control). Therefore, it is incumbent upon managers to make strategic decisions and allocate resources in the manner most consistent with the claims of the other stakeholder groups.

The unique role of managers suggests that they can be seen as the agents of other stakeholders; hence the term stakeholder–agency theory. It would be incorrect, however, to suggest that all the other groups of stakeholders are therefore principals in the sense implied by agency theory. In agency theory, principals hire agents to perform some service on their behalf. Stockholders and some customers apart, few stakeholders can be said to hire managers (in the case of employees the reverse is clearly true). Nevertheless, there is a parallel between the general class of stakeholder–agent relationships and the principal–agent relationships articulated by agency theory. Both stakeholder–agent and principal–agent relationships involve an implicit or explicit contract, the purose of which is to try and reconcile divergent interests. In addition, both relationships are policed by governance structures. Moreover, many of the concepts and much of the language of agency theory can be applied to stakeholder–agent relationships. All of this suggests that principal–agent relationships, as defined by agency theory, can be seen as a subset of the more general class of stakeholder–agent relationships.

### UNDERLYING ASSUMPTIONS

Our main assumptions concern the efficiency of the market mechanism. These assumptions have implications for the existence of power differentials between the parties to a contract. By a power differential, we mean a condition of *unequal dependence* between the parties to an exchange (Emerson, 1962; Pfeffer, 1981). That is, for two entities, $A$ and $B$, there is a power differential in $A$'s favour when $B$ depends upon $A$ more than $A$ depends upon $B$. (While this definition is suitable for the purposes of this article, a fuller discussion of the definitions of power can be found in the work of Lukes (1974) and Wrong (1988).

Agency theorists see the firm as surrounded by efficient markets that adjust quickly to new circumstances (Barney and Ouchi, 1986). They make the rather heroic assumption that markets are in or near an efficient equilibrium

(Fama, 1980; Fama and Jensen, 1983; Jensen, 1983). The efficient markets assumption implies that principals and agents have freedom of *entry* into and *exit* from contractual relationships. According to mainstream agency theorists, if an agent (principal) does *not* like the terms of a contract offered by a principal (agent) and/or the governance structures that police that contract, he/she can always seek a better alternative. If a shortage of agents (principals) results, the principals (agents) will be compelled by market forces to adopt more acceptable incentive mechanisms and/or governance structures. Of course, this argument ignores the obvious fact that a better alternative might not always be available (for a discussion of the implications, see Burt, 1983). Putting this criticism aside for a moment, however, the mainstream agency theory argument suggests that contracts between principals and agents, along with the governance structures that police those contracts, are determined by market forces. Thus, within the framework of mainstream agency theory, they must be seen as having efficiency properties.

This tidy logic breaks down if the efficient market assumption is dropped (Perrow, 1986; Putterman, 1984). If the markets that surround the firm are inefficient, as occurs when alternative contracting opportunities are limited, the existence of power differentials between principals and agents must be admitted. If agents are unable to exit from a contractual relationship without taking a substantial loss (because 'better alternatives' are not available), or if the supply of agents exceeds the demand for agents by principals, power shifts towards the principal. Similarly, if principals are unable to dismiss agents, or if there is a shortage of agents, power shifts towards the agents. This is important, since power differentials can materially affect both the content of principal–agent contracts and the structure of governance mechanisms policing those contracts.

Given this, it is important to state our assumptions with regard to market efficiency and equilibrium. There are two points of disagreement between stakeholder–agency theory and agency theory. The first concerns the speed with which markets adjust to new circumstances, while the second concerns the assumption of equilibrium. In contrast to agency theory, we view market adjustment processes as being characterized by friction (Williamson (1985) makes a similar point and argues that this friction results in transaction costs). Due to friction, once created disequilibrium conditions may persist for a prolonged period of time before an efficient equilibrium is re-established. The resulting disequilibrium conditions imply the existence of power differentials between the parties to an exchange.

*Sources of Friction*
Barriers to entry and exit constitute one source of friction (Porter, 1980). Barriers to entry and exit impede adjustment processes and may allow power differentials arising from an initial disequilibrium to persist in a market for significant periods of time. Open systems theory suggests a further source of friction. Managers and other stakeholders can to a degree shape or enact their environment (Pfeffer and Salancik, 1978; Weick, 1979). If a disequilibrium situation is perceived as being to their advantage, managers may be able to slow down the adjustment process by appropriate strategic investments (*e.g.*,

investments designed to increase entry barriers, by collusion, by predatory pricing, etc.).

Organizational inertia may be a further source of friction (Hannan and Freeman, 1984). As a result of disequilibrium, one party to a contract may be disadvantaged. Correcting the disadvantage may require the innovation of new incentive structures and/or monitoring and enforcement mechanisms. However, the ability of the disadvantaged party to innovate may be hindered by strong inertia forces. Due to inertia, it may be difficult to alter established routines and procedures for monitoring and enforcing management–stakeholder contracts so as to reflect new realities. Pressures such as sunk costs, political coalitions, the tendency to consider precedents as normative standards, and a simple lack of imagination all limit the degree to which incentive, monitoring, and enforcement structures respond quickly to new circumstances.

Having said this, in *the long run* we believe that market processes work to select out the most inefficient organizational forms. Despite barriers to entry and exit, attempts by managers to shape the environment to their advantage, and inertia, in the long run there are grounds for believing that the market system does achieve a *rough* balance between the efficient and inefficient (Alchain, 1950; Nelson and Winter, 1982). Although the adjustment process can be slowed down, it cannot be halted altogether. While we certainly do not agree that only the most efficient survive, we think it is highly probable that the most inefficient organizational forms lose ground and are eventually selected out.

## Equilibrium and Market Processes

If change was a rare event, the above arguments would imply that equilibrium situations in which the most inefficient organizational forms have been selected out are commonplace. However, one of the central features of the real world is that the only constant is change. Although market process work towards some kind of equilibrium, change constantly alters the direction in which that equilibrium is to be found. As argued by Schumpeter (1942), such change often occurs due to the process of creative destruction triggered by innovation. Moreover, Schumpeter suggested that innovation is itself a product of the competitive process. Thus, ongoing change may be a persistent endogenous feature of capitalism. Alternatively, change may be due to exogenous macro-environmental trends (demographics, social-political factors, macro-economic change, *etc.*). No matter how it arises, ongoing change creates a situation of permanent disequilibrium and hence, of persistent power differentials between stakeholders and managers. However, because of the random nature of change, power differentials are themselves unlikely to remain unidirectional. While change at one point in time may favour managers, change in a subsequent period may shift the balance of power towards other stakeholder groups.

Building on this perspective, we follow modern 'Austrian' economists in arguing that the focus of theoretical attention should be on *market processes* rather than equilibrium conditions in efficient markets, since the latter is little more than a convenient fiction (Kirzner, 1979; Knight, 1921; Littlechild and

Owen, 1980; Nelson and Winter, 1982; Schumpeter, 1942). While a drive towards efficiency may characterize the business system, in the sense that the most inefficient producers ultimately get selected out, we view short- and medium-term inefficiencies arising out of disequilibrium conditions as commonplace.[1] This suggests that power differentials arising out of disequilibrium conditions between managers and other stakeholders are an essential determinant of the nature of many stakeholder–agent contracts and the structures that police those contracts.

In sum, our view of market dynamics is fundamentally different from that which characterizes most of the agent–principal literature. We do not assume equilibrium, although we do assume that *market processes* work in such a manner that, in the long run, inefficient incentive structures and monitoring and enforcement mechanisms are selected out, while more efficient structures and mechanisms evolve to replace them. However, due to barriers to entry and exit, the ability of stakeholders and managers to enact their environment, and inertia, the adjustment process is plagued by significant frictions. By forcing attention onto adjustment processes under disequilibrium conditions, this perspective adds richness to the discussion of stakeholder–agency theory.

## DIVERGENT CLAIMS, UTILITY LOSS, AND CONTRACTING COSTS

The interests of principals and agents diverge primarily because these different groups have different utility functions. In turn, this can lead to direct conflict over the use to which resources are put (for example, see Jensen, 1986). Agency theory focuses on the divergence of interests between managers and stockholders. The argument can be traced back to the managerial discretion literature of the 1960s (Baumol, 1959; Marris, 1964; Williamson, 1964). This literature theorizes that stockholders are wealth maximizers, while managers maximize a utility function that includes remuneration, power, job security, and status as its central elements. The agency/ managerial literature postulates that satisfying the claims of stockholders involves maximizing the efficiency of the firm (Fama, 1980), while satisfying the claims of management requires increasing the size of the firm (remuneration, power, job security and status are argued to be a function of firm size). In turn, it has been argued that the desire to increase firm size results in a managerial preference for maximizing the growth rate of the firm, principally through diversification (Amihud and Lev, 1981; Aoki, 1984; Marris, 1964). Further, the discretion literature postulates a trade-off between growth maximization and efficiency maximization. Beyond a certain point, the greater investments in growth, the lower investments in maximizing efficiency. Thus, divergent claims give rise to an agency conflict between managers and stockholders.

Stakeholder–agency theory postulates that other stakeholder groups also place claims on the firm that, if satisfied, reduce the amount of resources that management can channel towards the pursuit of growth through diversification. Satisfying employee claims for higher wages, consumer claims for greater quality and/or lower prices, supplier claims for higher prices and

more stable ordering patterns, and the claims of local communities and the general public for lower pollution and an enhanced quality of life, all involve the use of resources that might otherwise be invested by managers in maximizing the growth rate of the firm. Thus, an agency conflict is inherent in the relationship between management and all other stockholders (for a transaction cost interpretation of this phenomenon see Williamson, 1985).

This is not to deny that to a degree the claims of stakeholders and managers also converge. For example, satisfying employee claims for higher wages and better working conditions may improve employee productivity and thus provide management with greater resources. Similarly, devoting resources to controlling pollution may result in local communities being more receptive to future proposals by management for expanding its operations. However, our contention is that beyond a certain point, this convergence of interest is replaced by divergence.

If uncorrected, the divergence between management and stakeholder preferences with regard to the way in which a firm allocates its resources will result in a failure of stakeholders to maximize their utility. The difference between the utility that stakeholders could achieve if management acted in stakeholders' best interests, and the utility that is achieved if management acts in its best interest, can be referred to as a *utility loss*. In the absence of incentive, monitoring, and enforcement structures that serve to align the interests of managers and stakeholders, utility loss may be substantial. The function of incentive, monitoring, and enforcement structures is to minimize utility loss by correcting for the divergence of interests between management and stakeholders.

The concept of utility loss leads to a somewhat broader definition of agency costs than that typically given in the agency literature. To distinguish this from the agency definition, from now on we shall talk in terms of *contracting costs*. These are defined as the reduction in utility that stakeholders bear by channelling resources to support incentive, monitoring, and enforcement structures, as opposed to using those resources directly to satisfy their utility function, plus any remaining or residual utility loss. For example, imagine that the maximum amount of utility that stakeholders can derive from a given relationship is 100 units, but that management preferences result in stakeholders only getting 60 units, resulting in a total utility loss of 40 units. If stakeholders devote resources equal to 10 units of utility to establishing incentive, monitoring, and enforcement mechanisms, they may increase the utility they derive from the relationship to 90 units, resulting in a net gain of 20 units. The remaining residual utility loss is 10 units. Thus, contracting costs are equal to the 10 units of utility that are sacrificed to support incentive, monitoring, and enforcement mechanisms, plus the residual utility loss of 10 units.

## INTEREST ALIGNMENT MECHANISMS

One way of minimizing the utility loss that arises from a divergence of interests involves introducing *ex-ante* interest alignment mechanisms into the

contracting scheme. In the agency literature, management and employee stock option plans are the most widely discussed of these mechanisms (Demsetz, 1983). Stock option plans serve to induce managers and employees to pay more attention to maximizing stockholder wealth, since that will simultaneously mazimize their own wealth. On a more general level, offering tax breaks for investments in pollution containment equipment is an example of how local communities and the general public (through their legislative agents) use incentives to try and align management interests with their own.

In addition, to gain access to their resources, stakeholders may demand that managers absorb *ex-ante bonding costs* in order to demonstrate their commitment to satisfying stakeholder interests. For example, consider the consumer contemplating entering an exchange relationship with a manufacturer of consumer durables. The purchase of durables presents consumers with a difficult agency problem. Consumer durables are purchased infrequently and involve large expenditure. In such circumstances, the consumer is vulnerable to opportunistic action on the part of management. Management may misrepresent the quality or durability of the product in an attempt to close the sale. The agency problem is solved by the *ex-ante* introduction of a warranty into the contracting scheme. This specifies management's obligation to correct defects or provide suitable compensation in the event of substandard quality. The warranty is a bonding mechanism that communicates to consumers a commitment on the part of management to a certain standard of quality.

More generally, a bonding mechanism is an example of the use of *credible commitments* (Williamson, 1985). Establishing a credible commitment requires that managers post a 'hostage' or bond forfeitable upon malperformance (Alchain and Woodward, 1988). The concept has been used to explain certain characteristics of a firm's relationships with its suppliers and consumers (although it is hardly limited to this context). For example, when a supplier has to make substantial investments in specialized assets in order to enter into trade with the firm, it is also exposing itself to the possibility of opportunistic abuse by management (Williamson, 1985). Once the supplier has made the investment, it is effectively 'locked in' to the relationship and cannot exit without reducing the value of those assets. Management may use this fact to go back on any *ex-ante* agreement and drive down the prices charged by the supplier. As insurance against this possibility, the supplier may demand a similar *ex-ante* investment in dedicated assets on the part of the firm. This locks both parties into a mutually dependent relationship in which power is symmetrically distributed. Examples include reciprocal trade agreements, most-favoured-buyer clauses, inflexible prices, posted prices, exclusive territories, franchise-specific investments, patent pools, and union shop agreements (Williamson, 1985). In all these cases, the underlying objective is to establish mutual dependency between managers and other stakeholder groups so that interests are more closely aligned. As a general rule, the use of credible commitments to bond managers and stakeholders will be greater the greater the investments in specialized assets required of either stakeholders or managers to support a given exchange relationship.

## MONITORING AND ENFORCEMENT MECHANISMS AND STRUCTURES

Interest alignment mechanisms apart, the contracts between stakeholders and managers are primarily implicit (Mitroff, 1983). Stakeholders supply the firm with resources on the implicit (tacit) understanding that their claims on the organization will be recognized. To ensure that this occurs, a number of institutional structures have evolved that serve the function of monitoring and enforcing the terms of implicit contracts. Agency theory generally refers to such institutions as governance structures. Our change of terminology reflects a broadening of emphasis. Specifically, we are concerned with more than just quasi-independent or third party governance (such as the board of directors, the market for corporate control, or the legal superstructure of society). We are also concerned with institutions that have evolved to represent and further the interests of a given set of stakeholders (such as labour unions and consumer unions) precisely because such institutions have utility loss-minimizing properties. Thus, the term institutional structures subsumes the term governance structures.

*Monitoring Structures*
An information asymmetry exists between managers and stakeholders. As insiders, managers are in a position to filter or distort the information that they release to other stakeholders. Management control over critical information complicates the agency problem. It makes it difficult for stakeholders to identify if management is acting in their interests. The obvious response is for stakeholders to gather more information about management activities. However, while individual stakeholders can and do undertake their own monitoring of management performance, the costs of gathering and analysing additional information may be prohibitive.

This is particularly likely when stakeholders are diffused. Diffusion refers to a situation where a stakeholder group contains many individuals or entities, no one of which has command over a significant proportion of the group's total resources. In such circumstances, *ceteris paribus*, no one individual or entity may be able to finance the extensive information-gathering and analysis necessary to reduce significantly the information asymmetry between managers and stakeholders. In turn, this gives managers greater discretionary control over the use to which the firm's resources are put, increasing the residual loss that stakeholders have to bear.

The response to the monitoring problem has been the evolution of a wide range of institutional structures that serve to economize upon the costs of information-gathering and analysis. Some of these structures are enshrined in legislation (*e.g.*, the requirement that public companies publish consolidated annual accounts). Other institutions have evolved in an attempt to exploit the profit opportunities of gathering, analysing, and then selling information to stakeholders (*e.g.*, stock analysts' services, consumer reports, *etc.*). Still others have arisen as non-profit organizations that exist in part to monitor the degree to which managers act in the best interests of certain stakeholder groups (*e.g.*, Consumer Watch, Infact, labour unions). The common theme found in all of these structures is their ability to achieve economies of scale in information-

gathering and analysis, primarily through the employment of specialists. The consequence of such devices is a reduction in utility loss.

### Enforcement Mechanisms and Structures

The function of enforcement is one of deterrence. Enforcement mechanisms are articulated by stakeholders prior to any resource exchange in an attempt to deter management from maximizing its utility at the expense of stakeholders. The success of enforcement mechanisms depends upon their credibility (Schelling, 1960), and those lacking credibility will be ignored by management. In such circumstances, any attempt to put enforcement mechanisms into effect will involve costs that outweigh the benefits of reducing the utility loss from management opportunism. In short, mechanisms that are not effective deterrents will fail (as do laws that are commonly ignored by the general population).

*Law as a deterrent.* Establishing credible deterrents in the context of stakeholder–management relationships requires enforcement mechanisms that are supported by a broad consensus of stakeholders, and which are effectively communicated to management *ex-ante*. Certain legal penalties have this character (laws against insider trading, antitrust regulations, pollution regulations, *etc.*). Indeed, it can be argued that much of the structure of law relating to business activity in society reflects critical points of conflict in stakeholder–agent relationships. That is, legislators, as representatives of certain stakeholder interests, have enacted into law enforcement mechanisms that, because they are credible deterrents, serve to economize on utility loss.

*Exit as a deterrent.* The legal approach to resolving principal–agent conflicts constitutes only one way of establishing a credible threat. A more general approach involves the establishment of a credible threat to withhold resources from the firm if management fails to serve stakeholder interests. That is, to threaten exit from the exchange relationship (Hirschman, 1970). Such threats may be more effective than legal penalties. Only in rare situations are legal penalties likely to jeopardize the survival of the firm. Indeed, many firms view such penalties as a 'normal cost of doing business'. In contrast, by denying the firm access to critical resources, stakeholders can threaten its very survival (Pfeffer and Salancik, 1978).

In a sense, the threat of exit is an underlying theme of many stakeholder–management relationships. For example, if dissatisfied with product quality, consumers can always take their business elsewhere. Similarly, if dissatisfied with a firm's performance, stockholders can always sell their stock. Thus, credible threats to exit can be enacted through the market mechanism. However, market action suffers from a number of weaknesses. First, there is a co-ordination problem among diffused stakeholders that in certain circumstances makes collective action problematic. Exit may not be a very effective deterrent if members of a stakeholder group are unable to act in unison to impose demands on management. For example, while employees may be unhappy about working conditions, individual complaints or threats of exit may do little to persuade management to improve conditions, particularly if

there is a ready supply of replacement labour. Similarly, while consumers as individuals may disapprove of the pollution implications of a given product (*e.g.* auto exhaust, plastic containers, air conditioning fluid) the threat of individual exit may be futile, particularly if no cost-effective alternative exists. Thus, consumers may continue to purchase the products, even though as individuals they are unhappy about the implications of doing so, and would prefer the firm to devote resources to developing less harmful alternatives.

The institutional response to the problem of achieving collective action among diffused stakeholders has been the evolution of a number of structures that perform the co-ordination function. Examples include labour unions, consumer unions, and special-interest groups. By providing centralized direction, these structures economize upon the costs of co-ordination and establish the credibility of the exit mechanism. Thus, labour unions may initiate a strike if management fails to meet their demands for better working conditions. Similarly, special-interest groups may initiate a consumer boycott if the firm continues to produce products that they consider to be harmful (*e.g.*, Infact's consumer-led boycott of Nestle's was designed to halt the company's questionable infant formula marketing practices in Third World countries).

A more intractable weakness of market action is that it may lack effectiveness in those situations where stakeholders are 'locked in' to an exchange relationship by specific asset investments. Suppliers, customers, employees, or communities who have invested in specialized assets in order to enter into an exchange relationship with the firm may not be able to exit without incurring substantial exit costs. The exit costs consist of the reduced rents from specialized assets that can be earned in their next best application. Other things being equal, such barriers to exit reduce the credibility of any threat to exit as a contract enforcement mechanism. This is serious given that actors who make specific asset investments in the firm are by definition among the most important of its stakeholders (their future is most closely aligned to that of the firm).

However, certain bonding mechanisms have the additional character of increasing the credibility of the threat to exit among stakeholders who have invested in firm-specific assets. For example, a union shop agreement can be viewed as a bonding mechanism by which management agrees to hire only union labour as a means of safeguarding employees investments' in firm-specific human capital. This bonding mechanism limits management's ability to abrogate any previously agreed labour contract. If they do, they face the possibility of a strike (exit), the threat of which is made credible by the inability to hire non-union labour. Notwithstanding such examples, however, the threat of exit may be limited in such circumstances, in which case stakeholders may have to resort to voice as an enforcement mechanism (Hirschman, 1970).

*Voice as a deterrent.* In certain circumstances voice may be the most effective enforcement mechanism. Voice is often the least costly mechanism to adopt. Newsworthy publicity comes cheap, yet it can severely damage managerial reputations and the intrinsic value of a manager's human capital. To be effective, however, voice must be articulated by interest groups that have a

legitimate claim to represent stakeholder interests. Again, certain institutional structures such as labour unions, consumer unions, and special-interest groups arguably have this characteristic. This reinforces our earlier conclusion that interest groups can be viewed as *institutional structures that have evolved to economize on contracting costs*.

## STATIC EQUILIBRIUM

In our view, due to the pervasive nature of change, much of the business system is in a state of almost permanent disequilibrium. Despite this, there is value in discussing what we would expect to find if the business system were ever to achieve equilibrium. Although this is something of an abstract and teleological exercise, such a discussion tells us something about the end towards which dynamic processes propel the system. Here we discuss the factors determining the complexity of the institutional structures that we would expect to find in an equilibrium situation; later we focus on disequilibrium.

### A Static Model

If equilibrium were ever reached, institutional structures would display efficiency properties. Specifically, stakeholders would increase the complexity of institutional structures up to that point where the marginal benefits of doing so (in terms of a reduction in utility loss) were equivalent to the marginal costs of maintaining those structures (in terms of the utility that has to be sacrificed to support them). Given this, it is probable that in equilibrium managers still retain some (diminished) discretionary control over the use to which the firm's resources are put. The argument is explained with reference to figure 1.

The horizontal axis of figure 1 measures the complexity of *existing* institutional structures. The least complex structure is that of the market mechanism. More complex structures involve increasingly extensive monitoring and enforcement mechanisms. Thus, consumer watchdogs such as Ralph Nader's Consumer Watch, or the development of labour unions, can be seen as adding complexity to the institutional structures that police the management–stakeholder interface. The vertical axis measures units of utility.

A positive relationship between the complexity of available institutional structures and the costs of those structures (in terms of the utility that has to be sacrificed to support them) can be postulated. If working efficiently, the market system, because it is a decentralized mechanism, imposes the lowest costs on stakeholders. More complex structures impose additional costs. For example, ultimately consumers underwrite Consumer Watch through donations. Employees underwrite labour unions through subscriptions. Similarly, if stakeholder pressures result in certain regulations being enacted into law, ultimately stakeholders, as taxpayers, underwrite the commensurate increase in legal apparatus.

However, due to the benefits of specialization it seems likely that economies of scale in information-gathering and analysis exist. Thus, initially the costs

Figure 1

(in terms of utility) of maintaining institutional structures will increase at a decreasing rate with increasing complexity. This is illustrated in figure 1 where the cost function increases at a decreasing rate up to the point of inflection *a*. Past *a* diminishing returns to specialization are likely to set in and costs will increase at an increasing rate.

The benefits to stakeholders of maintaining institutional structures can be measured in terms of the *reduction in utility loss* that such structures achieve. The benefit function in figure 1 is shown to increase at a decreasing rate, symbolizing decreasng returns to increasingly complex structures; that is, increasing management resistance to reductions in their discretionary control over the use to which resources are put. Eventually, the function will approach the line *bb'*, where 0–*b* symbolizes the total utility loss arising from an *ex-ante* divergence of interests.

The equilibrium condition in figure 1 involves the stakeholders devoting resources to increasingly complex institutional structures up to that point where the marginal benefits of such expenditures are equivalent to the marginal costs. It is worthwhile for stakeholders to bear the costs of establishing and running institutional structures of 0–*C*1 complexity.[2] Note, this equilibrium point involves a reduction in total utility loss of 0–*c*. The remaining utility loss is equivalent to *c*–*b*. Thus, *c*–*b* represents the resources still under the discretionary control of management once the claims of stakeholders have been satisfied. The logic of our earlier arguments suggests that these resources will be devoted to investments in maximizing the growth rate of the firm.

Another way of viewing *c*–*b* is as a measure of the incentive stakeholders have to develop more efficient institutional mechanisms. More precisely, the

*gross* returns to innovation are equivalent to the discounted present value of $\Sigma c_t - b_t e$ where the subscript $t$ refers to successive time periods. For $C1$ to represent a true equilibrium, the perceived returns to innovation must be equivalent to the perceived costs of innovation. The costs of innovation refer to the costs of overcoming resistance to change (in terms of the utility that must be sacrificed) and imposing new institutional structures upon the implicit or explicit contract. An example of these costs might be the costs in terms of both money and emotion to employees of supporting a strike to get their labour contract with management renegotiated. If the perceived returns to innovation are greater than the foreseeable costs, it will pay stakeholders to devote resources to the development of more efficient institutional structures. The implications of this point are developed later.

### Extensions

A shortcoming of this model is that it glosses over the problems created by the conflicting claims of different stakeholder groups. Obviously, the claims of different groups may conflict (*e.g.*, stockholder demands for greater dividends conflict with employee demands for higher wages). However, on a more general level, each group can be seen as having a stake in the continued existence of the firm.[3] Where opinions differ between stakeholder groups is on how the firm's resources should be allocated between investments, and the most desirable time pattern of organizational rent streams. If the different stakeholder groups engage in open conflict over this issue, the net effect may well be to damage the firm and all involved with it (as when employees go on strike or consumers boycott its products). Thus, different stakeholder groups have an incentive to *co-operate*, rather than incur the costs of open conflict (for a theoretical discussion of this see Aoki, 1984).

An equilibrium solution to this type of problem can be found in the literature on co-operative game theory. Although beyond the scope of this article, it should be noted that it is possible to model what has been referred to as an 'organizational equilibrium' (Aoki, 1984). This is a state in which no one group of stakeholders can increase its utility without risking a higher expected loss of utility owing to the possible withdrawal of co-operation by the other stakeholders. A rational stakeholder would not disturb such a state by making a demand for greater control over how the firm's resources are invested (for a theoretical proof of this argument see Aoki, 1984).

With reference to institutional structures, the implication of such an organizational equilibrium is that each stakeholder group will adopt increasingly complex structures up to the point that is consistent with the co-operative solution. That is, no one group will attempt to establish additional institutional structures if doing so would upset the organizational equilibrium and precipitate open conflict between stakeholders.

Management's role in this process is one of an interest mediator. Management is assigned the difficult task of balancing conflicting demand so as to achieve a co-operative solution. Management is hardly a passive player, however. Management can be viewed as trying to expand its bargaining position with respect to different stakeholder groups. Under the restrictive conditions of neoclassical equilibrium, such an exercise would be fruitless.

However, an Austrian perspective of the market process leads to a very different conclusion.

## POWER DIFFERENTIALS AND MARKET PROCESS

It was argued earlier that due to the pervasiveness of change, extensive disequilibrium is the norm. Moreover, although we view markets as being ultimately efficient, we theorized that the adjustment process is characterized by considerable friction due to inertia, the ability of managers and stakeholders to slow down adjustment by their strategic investments, and entry and exit barriers. The implication of prolonged disequilibrium is that in practice, power differentials arising from a condition of dependency between principals and agents are commonplace and may persist for some time. The advantaged party may use such differentials to further entrench its position and modify institutional structures to its advantage.

Of course, power differentials do not always work to management's advantage. For example, labour shortages arising from unanticipated macro-environmental change will increase the bargaining power of employees relative to managers, enabling them to impose tighter constraints on managers (*e.g.*, to demand higher wages, better working conditions, more extensive grievance procedures, and employee directors). Often, however, power differentials will be in management's favour. Moreover, by virtue of their position at the nexus of the implicit and explicit contracts that constitute the firm, and because of their control over the decision-making apparatus of the firm, managers may be better positioned to exploit power differentials than individual groups of stakeholders. Thus, in the remainder of this section we will focus on management strategies for establishing and/or exploiting power differentials, the implications of power differentials for institutional structures, and stakeholder responses to management actions.

### Establishing and Exploiting Power Differentials

Starting with the convenient fiction that in the 'beginning there was an efficient equilibrium', disequilibrium can be seen as either the product of a firm's own innovative efforts, or the result of an exogenous shock. However created, management may try to take advantage of the resulting turbulence and uncertainty to engineer a situation in which the firm's stakeholders are more dependent on management than management is upon them. This involves undertaking strategic actions that reduce the concentration of stakeholder power and/or increase the concentration of management power.

The concentration of stakeholder power can be reduced by strategies designed to *diffuse* the control over critical resources exercised by stakeholder groups. For example, with reference to stockholders, *targeted* stock buybacks along with new stock issues may be used to reduce ownership concentration and increase shareholder dispersion. Dispersion makes it more difficult for stockholders to monitor and enforce their implicit contract with management (Berle and Means, 1932). In a similar vein, management may diffuse supplier power by developing alternative sources of supply (assuming that alternatives

are available). Management may reduce customer power by building a more diverse customer base through product and market diversification. Management may limit the power of local communities and the general public by both national and multinational diversification. And finally, it has been argued that the way in which management has organized production in the workplace and has exercised control through bureaucratic mechanisms has significantly reduced the power of employees to oppose management policies (Braverman, 1974; Clawson, 1980; Edwards, 1979).

Increasing the concentration of management power requires strategies that increase the amount of resources under management control. These include horizontal mergers and acquisitions to increase concentration within an industry, vertical integration to gain power over suppliers and customers, and co-operative agreements between the managers of different firms including joint ventures, interlocking directorates, purchasing alliances, and price leadership agreements (Pfeffer and Salancik, 1978). The common theme underlying these strategies is that they restrict the choice set of stakeholders, thereby altering the configuration of resource dependencies. For example, horizontal acquisitions increase the buying power of the firm by limiting the number of independent customers to whom suppliers can sell.

All of these strategies are undertaken to increase management power rather than maximize efficiency. Their ultimate objective is to loosen the constraints imposed by stakeholders and give management greater discretionary control over the firm's resources. Without a commensurate increase in productive efficiency, the additional bureaucratic costs of running an expanded organization or of achieving intra-organizational co-ordination imply that declining efficiency will be one result of such strategies. Thus, in an efficient market, firms that pursue such strategies will be selected out by the competitive mechanism. However, the view of competitive dynamics advocated here suggests that disequilibrium gives managers the opportunity to build such power differentials.

Of course, it is possible that the ability of managers to pursue strategies that increase management power over one group of stakeholders may be limited by the constraints imposed by other stakeholder groups. Most significantly, the board of directors (as the representative of stockholders), is in theory well positioned to limit managerial actions that it perceives as being contrary to stockholder interests (Fama and Jensen, 1983). Thus, for example, management attempts to reduce customer power by building a more diverse customer base through diversification may be blocked by the board, precisely because the board might regard such diversification as being an inefficient use of stockholders' funds.

Whether the board can impose such constraints in practice is the subject of some debate. Contrary to the argument made by Mace (1971) and others that most boards do little more than rubber-stamp management decisions, Mizruchi (1983) has presented strong arguments in support of the proposition that board control of management actions in public corporations is still possible. More recently, Lorsch and MacIver (1989) present case study evidence which suggests that among selected United States corporations, boards are increasingly exercising their control over top management teams.

On the other hand, there is also evidence which suggests that board control over top management is still relatively weak. For example, Jensen and Murphy (1990), after finding only a very weak relationship between CEO pay and firm performance, concluded that most boards may lack the power to impose stockholder objectives on management. Similarly, Burrough and Helyar (1990) have described in detail how one CEO, Ross Johnson of RJR Nabisco, handed out lucrative consulting contracts to outside directors in a successful attempt to keep them from criticizing management policies that were clearly inconsistent with stockholders' best interests. The issue of how strong boards actually are, therefore, remains an open one.

However, one factor which suggests that tighter control over management actions may become the rule rather than the exception has been the dramatic rise of financial institutions as major providers of capital. On both sides of the Atlantic pension funds, insurance companies, mutual funds, and investment banks have rapidly been replacing individuals as the main stockholders in public corporations. For example, in the United States, Hanson and Hill (1991) present evidence which suggests that among Fortune 500 companies the percentage of common stock held by institutions increased from 24 percent to 50 percent between 1977 and 1986. The growing concentration of stockholding in the hands of a relatively few institutions is resulting in the evolution of a stock market that bears little resemblance to the fragmented and dispersed market described by Berle and Means (1932). Instead, the resulting concentration of stockholdings means that financial institutions are increasingly able to exert direct influence over management actions, either through (a) the threat to sell their holdings; (b) the threat to fight proxy votes more aggressively; or (c) by using their voting power to elect their own nominees to the board of directors.

For example, in 1987 a group of financial institutions with major holdings in General Motors was able to pressure GM management into adopting a bonus pay system for GM executives that was based upon stock price performance (prior to that time, bonuses had been awarded automatically, irrespective of the company's performance). The institutions did this by threatening to introduce a resolution at the next stockholders' meeting that would be critical of management unless the company changed its bonus pay policies (Nussbaum and Dobrzynski, 1987). More generally, Mintz and Schwartz (1985) argue for and present evidence which supports the view that financial institutions play a key role in the control of large firms. Similarly, Scott (1979) concludes that large firms and major banks 'confront one another as equals, each being constrained by its controlling constellation of interests' and that 'banks are able to exercise considerable influence over the policies of major industrial corporations and so can affect what happens in companies where they have no direct power' (p. 175).

It should be pointed out, however, that to a large degree management and major financial institutions share the same agenda. Although there will undoubtedly be conflict between them, it is reasonable to suppose that in many cases management actions designed to weaken the power of certain stakeholder groups (e.g., employees, suppliers, or customers), will be congruent with the interests of major financial institutions so long as they

increase the profitability of the corporation. Thus, while important, the potential for conflict between managers and financial institutions should not be overstated.

*The Implications of Power Differentials*

Power differentials created by the strategies detailed above limit the ability of stakeholders to enforce implicit or explicit contracts. Diffusion of stakeholder power makes co-ordination between individual stakeholders more problematic and costly, thereby reducing the ability of stakeholders to act collectively. In turn, this limits the effectiveness of voice and exit as enforcement mechanisms. It is more difficult for stakeholders to establish a credible threat when power is diffused among many individuals and collective action is difficult to achieve. Similarly, the concentration of management power reduces the choice set of stakeholders, again limiting the effectiveness of exit and voice as enforcement mechanisms.

Stakeholder diffusion also makes monitoring more difficult. Less powerful stakeholders are less able to demand that management make itself accountable. They are less able to use the implied threat to exit or exercise voice as a means of gaining access to insider information or demanding that management regularly provides them with information concerning its activities. Moreover, the pursuit of diversification strategies by the firm obscures data relating to the efficiency of individual divisions (firms only have to publish consolidated accounts). This exacerbates the information asymmetry between management and stakeholders, making monitoring more problematic.

Managers may also take advantage of power differentials unilaterally to rewrite the terms of the implicit or explicit contract between managers and stakeholders. Thus, managers may take advantage of power differentials to revoke warranties, retract hostages posted as bonds, or retract other credible commitments such reciprocal purchasing agreements, posted prices, or union shops. Similarly, management may take advantage of a temporary power differential over its employees to rewrite employment contracts. In all of these cases, the effect of power differentials is to reduce the effectiveness of existing institutional structures and to increase the residual loss that must be born by stakeholders.

*Stakeholder Responses*

Stakeholder responses to the creation of power differentials can be analysed by way of figure 2. This shows the marginal benefit and marginal cost curves underlying figure 1. We start the analysis by accepting the convenient fiction of an initial equilibrium solution involving institutional structures of $C1$ complexity, a reduction in utility loss of $0-c$, and a remaining utility loss of $c-b$. The effect of a successful attempt by management to create a power differential will be to reduce the gradient of the benefit function, and hence shift the marginal benefit function down from $MB_1$ to $MB_2$. In other words, power differentials limit the effectiveness of existing institutional structures and result in a reduction in the utility loss that can be achieved by stakeholders at each level of institutional complexity. The new equilibrium solution implied by this shift is to be found at $C2$. Thus, comparative statics suggest

Figure 2

that when faced with an adverse power differential it pays stakeholders to reduce institutional complexity to C2 and accept an increased residual loss of $d-b$.

However, in a dynamic sense the existence of $d-b$ can be seen as providing an incentive to stakeholders to find new ways of economizing on contracting costs (to develop new institutional structures). As noted earlier, for C1 to be an equilibrium position, the perceived gross return gained from the innovation of more efficient structures must be equivalent to the perceived costs of such innovation. Given this, the power shift has created an incentive to innovate equivalent to the discounted present value of $\Sigma c_t - d_t$. Our thesis is that the existence of such an incentive following the emergence of a power differential has driven much of the historical evolution of institutional structures. The evolution of labour unions, consumer unions, special-interest groups, incentive mechanisms and credible commitments, corporate regulation, and so on, can be traced back to such incentives.

For example, the development of the factory system in nineteenth-century England led to the decline of the 'putting-out' system of subcontracting and upset the balance of power that existed between those who made products to those who managed the production process, with management benefiting (Landes, 1966). Those who made products now became 'employees' and were at a power disadvantage vis-à-vis those who managed the process. Traditionally, craft guilds had governed the implicit contract between 'managers' and 'subcontractors'. One consequence of the shift to a factory system was that craft guilds lost their effectiveness as institutional structures, and declined dramatically in influence (Landes, 1966). Thus, the decline in the marginal benefit curve from $MB_1$ to $MB_2$ (i.e., existing institutional structures were no longer effective). However, this shift increased the incentive that those who

261

made products (employees) had to develop new and more effective govern-
ance mechanisms. The response was the development of labour unions with
the objective of re-establishing a condition of mutual dependence between
those who made products and those who managed the manufacturing pro-
cess. The effectiveness of unions was based on their ability to economize on
co-ordination costs between diffused stakeholders and re-establish exit as a
credible threat.

Of course, such adjustments are anything but smooth and will be resisted
by the advantaged party. Indeed, there is a long history of management
resistance to the development of union power following the introduction of the
factory system. Moreover, the conflict between management and stakeholders
following the development of a power differential can be expected to spill over
into the explicitly political arena. That is, both parties can be expected to try
to use the power of law to further their interests. This is hardly surprising
given that many institutional structures either have a legal component or are
supported by the law, but it does give us a way of explaining the selective use
of Political Action Committee (PAC) money, along with more general lob-
bying by corporate trade associations and public interest groups. Specifically,
at any point in time such monies and lobbying will be devoted to ongoing
management–stakeholder conficts, the amount of activity and money being
roughly proportional to the size of the perceived power differential and the
anticipated gains from either changing the system or maintaining the *status
quo*.

Finally, it is important to remember that power differentials work both
ways. Although we have concentrated on the benefits enjoyed by manage-
ment from their control over the decision-making apparatus of the firm, and
although this control probably does give management an inbuilt advantage,
management may be put on the defensive by increases in stakeholder power
(as may be occurring *vis-à-vis* stockholders due to the increase in the amount
of stock held by financial institutions). As with management power, increases
in stakeholder power have their genesis in disequilibrium conditions created
either by exogenous shocks, or by innovations in the way that stakeholders do
business.

## CONCLUSION

The objectives of this article were ambitious. Taking agency and stakeholder
perspectives of the firm as our starting point, we have attempted to construct
a paradigm that explains certain aspects of the strategic behaviour of the firm,
the structure of incentive alignment mechanisms, and the institutional forms
that have evolved to police the implicit and explicit contracts between
managers and stakeholders. In doing so, we have drawn on the literatures of
business and society, economics, finance, and organizational theory.

The resultant paradigm, stakeholder–agency theory, can be viewed as a
modification of agency theory to accommodate theories of power including
resource dependence theories of organizations. Hitherto, these theories have
been seen as offering mutually exclusive interpretations of organizational

phenomena (Perrow, 1986). While agency theory assumes efficient markets and rejects the idea of power differentials between managers and stakeholders, resource dependency theory (*e.g.*, Pfeffer and Salancik, 1978) implicity assumed inefficient markets which allow for the existence of unequal resource dependencies (power differentials) between managers and stakeholders. The adoption of an 'Austrian' perspective on market processes allows us to treat notions of power and efficiency within the framework of the same model.

Following the theme of Austrian economics, we accept that markets are efficient. However, the existence of short-run disequilibrium arising from exogenous and endogenous change has been argued to give rise to temporary power differentials between managers and stakeholders. Some of the strategies pursued by managers with respect to stakeholders can be seen as an attempt to exploit and entrench these power differentials. In turn, the evolution of new incentive structures and institutional mechanisms for monitoring and enforcing the contractual relationships between managers and stakeholders can be seen as long-run market-generated responses to disequilibrium conditions and unequal resource dependencies.

Our contention is that joining together notions of power and efficiency within the same framework substantially increases the predictive power of the paradigm when compared to earlier 'theories of the firm'. Unlike earlier theories, the paradigm explicitly focuses on the causes of conflict between managers and stakeholders following the emergence of disequilibrium conditions. Stakeholder–agency theory also points the way towards a theory of the adjustment mechanisms that realign management and stakeholder interests following disruption.

## NOTES

* Our thanks to Peter Mills, Tom Thomas and an anonymous referee for their helpful comments on an earlier draft of this manuscript.

[1] There is in fact a large body of empirical evidence suggesting this is indeed the case. See Branch (1980), Bronzen (1970) and Jacobsen (1988).

[2] $C1$ is where the gradients of the two curves are equivalent: *i.e.* where $d$ (benefits)/$d$ (complexity) $= d$ (costs)/$d$ (complexity).

[3] There are exceptions to this, however. For example, stockholders often stand to gain large profits if they sell out to corporate raiders. At this juncture, the continued fate of the firm is of little interest to them.

## REFERENCES

ALCHIAN, A. A. (1950). 'Uncertainty, evolution theory'. *Journal of Political Economy*, **58**, 211–21.

ALCHIAN, A. A. and WOODWARD, S. (1988). 'The firm is dead, long live the firm: a review of Williamson's *The Economic Institutions of Capitalism*'. *Journal of Economic Literature*, **26**, 65–85.

AMIHUD, Y. and LEV, B. (1981). 'Risk reduction as a managerial motive for conglomerate mergers'. *Bell Journal of Economics*, **12**, 393–406.

AOKI, M. (1984). *The Cooperative Game Theory of the Firm*. Oxford: Oxford University Press.

BARNEY, J. B. and OUCHI, W. G. (1986). *Organizational Economics*. San Francisco: Jossey-Bass.

BAUMOL, W. J. (1959). *Business Behavior, Value and Growth*. New York: Macmillan.

BERLE, A. A. and MEANS, G. C. (1932). *The Modern Corporation and Private Property*. New York: Macmillan.

BRANCH, B. (1980). 'The laws of the marketplace and ROI dynamics'. *Financial Management*, **9**, 58–65.

BRAVERMAN, H. (1974). *Labor and Monopoly Capital*. New York and London: Monthly Review Press.

BRONZEN, Y. (1970). 'The antitrust task force deconcentration recommendation'. *Journal of Law and Economics*, **13**, 279–92.

BURROUGH, B. and HELYAR, J. (1990). *Barbarians at the Gate: The Fall of RJR Nabisco*. New York: Harper & Row.

BURT, R. S. (1983). *Corporate Profits and Cooptation*. New York: Academic Press.

CLAWSON, D. (1980). *Bureaucracy and the Labor Process*. New York and London: Monthly Review Press.

DEMSETZ, H. (1983). 'The structure of corporate ownership and the theory of the firm'. *Journal of Law and Economics*, **26**, 375–89.

EDWARDS, R. (1979). *Contested Terrain*. New York: Basic Books.

EISENHARDT, K. M. (1985). 'Control: organizational and economic approaches'. *Management Science*, **31**, 134–49.

EISENHARDT, K. M. (1988). 'Agency and institutional explanations of compensation in retail sales'. *Academy of Management Journal*, **31**, 488–511.

EISENHARDT, K. M. (1989). 'Agency theory: an assessment and review'. *Academy of Management Review*, **14**, 57–74.

EMERSON, R. M. (1962). 'Power dependence relations'. *American Sociological Review*, **27**, 31–41.

FAMA, E. F. (1980). 'Agency problems and the theory of the firm'. *Journal of Political Economy*, **88**, 375–90.

FAMA, E. F. and JENSEN, M. C. (1983). 'Separation of ownership and control'. *Journal of Law and Economics*, **26**, 301–26.

FREEMAN, E. (1984). *Strategic Management: A Stakeholder Approach*. Boston: Pitman Press.

HANNAN, M. T. and FREEMAN, J. (1984). 'Structural inertia and organizational change'. *American Sociological Review*, **49**, 149–64.

HANSON, G. and HILL, C. W. L. (1991). 'Are institutional investors myopic? A pooled time series analysis in four industries'. *Strategic Management Journal*, **12**, 1–16.

HIRSCHMAN, A. (1970). *Exit, Voice and Loyalty*. Cambridge, Mass.: Harvard University Press.

JACOBSEN, R. (1988). 'The persistence of abnormal returns'. *Strategic Management Journal*, **9**, 415–30.

JENSEN, M. C. (1983). 'Organization theory and methodology'. *Accounting Review*, **50**, 319–39.

JENSEN, M. C. (1986). 'Agency costs of free cash flow, corporate finance and takeovers'. *American Economic Review*, **76**, 323–9.

JENSEN, M. C. and MECKLING, W. H. (1976). 'Theory of the firm: managerial behavior, agency costs, and ownership structure'. *Journal of Financial Economics*, **3**, 305–60.

JENSEN, M. C. and MURPHY, K. (1990). 'Performance pay and top management incentives'. *Journal of Political Economy*, **98**, 225–64.

KIRZNER, I. M. (1979). *Perception, Opportunity, and Profit*. Chicago: University of Chicago Press.

KNIGHT, F. (1921). *Risk, Uncertainty, and Profit*. Boston: Houghton-Mifflin.

KOSNIK, R. (1987). 'Greenmail: a study of board performance in corporate governance'. *Administrative Science Quarterly*, **32**, 163–85.

LANDES, D. S. (1966). *The Rise of Capitalism*. New York: Macmillan.

LITTLECHILD, S. C. and OWEN, G. (1980). 'An Austrian model of the entrepreneurial market process'. *Journal of Economic Theory*, **23**, 361–79.

LORSCH, J. W. and MACIVER, E. (1989). *Pawns or Potentates: The Reality of America's Corporate Boards*. Boston, Mass.: Harvard Business School Press.

LUKES, S. (1974). *Power, A Radical View*. London: Macmillan.

MACE, M. L. (1971). *Directors: Myth and Reality*. Cambridge, Mass.: Harvard Business School Press.

MARCH, J. G. and SIMON, H. A. (1958). *Organizations*. New York: Wiley.

MARRIS, R. (1964). *The Economic Theory of 'Managerial' Capitalism*. London: Macmillan.

MINTZ, B. and SCHWARTZ, M. (1985). *The Power and Structure of American Business*. Chicago: University of Chicago Press.

MITROFF, I. I. (1983). *Stakeholders of the Organizational Mind*. San Francisco: Jossey-Bass.

MIZRUCHI, M. S. (1983). 'Who controls whom? An examination of the relation between management and boards of directors in large American corporations'. *Academy of Management Review*, **8**, 426–35.

NELSON, R. R. and WINTER, S. G. (1982). *An Evolutionary Theory of Economic Change*. Cambridge, Mass.: Harvard University Press.

NUSSBAUM, B. and DOBRZYNSKI, J. (1987). 'The battle for corporate control'. *Business Week*, 18 May, 102–9.

PEARCE, J. A. (1982). 'The company mission as a strategic tool'. *Sloan Management Review*, Spring, 15–24.

PERROW, C. (1986). *Complex Organizations: A Critical Essay*. New York: Random House.

PFEFFER, J. (1981). *Power in Organizations*. Marshfield, Mass.: Pitman Publishing.

PFEFFER, J. and SALANCIK, G. R. (1978). *The External Control of Organizations: A Resource Dependence Perspective*. New York: Harper & Row.

PORTER, M. A. (1980). *Competitive Strategy*. New York: Free Press.

PUTTERMAN, L. (1984). 'On some recent explanations of why capital hires labor'. *Economic Inquiry*, **33**, 171–87.

ROSS, S. (1973). 'The economic theory of agency: the principal's problem'. *American Economic Review*, **63**, 134–9.

SCHELLING, T. C. (1960). *The Strategy of Conflict*. Cambridge, Mass.: Harvard University Press.

SCHUMPETER, J. A. (1942). *Capitalism, Socialism, and Democracy*. New York: Harper Brothers.

SCOTT, J. (1979). *Corporations, Classes and Capitalism*. London: Hutchinson.

WEICK, K. E. (1979). *The Social Psychology of Organizing*. Reading, Mass.: Addison-Wesley.

WILLIAMSON, O. E. (1964). *The Economics of Discretionary Behavior: Managerial Objectives in a Theory of the Firm*. Englewood Cliffs, N.J.: Prentice-Hall.

WILLIAMSON, O. E. (1984). 'Corporate governance'. *Yale Law Review*, **93**, 1197–230.

WILLIAMSON, O. E. (1985). *The Economic Institutions of Capitalism*. New York: Free Press.

WRONG, D. (1988). *Power: Its Forms, Bases, and Uses*. Chicago: University of Chicago Press.

# The agency problem, agency cost and proposed solutions thereto: A South African perspective

**J H Hall**

**Abstract**

The development and growth of listed firms during the past few decades has caused an ever-widening gap between ownership and management. The agency theory addresses this relationship between owners (shareholders) and the custodians of their wealth, that is the management of a firm. If management's goals differ from those of the firm, an agency problem arises and the owners have to incur agency cost to overcome this problem.

Besides discussing the theoretical principles underlying the above issues, an empirical investigation was undertaken, using questionnaires completed by firms listed on the Johannesburg Stock Exchange.

It appears from the responses received as if the agency problem does exist in a significant number of companies. Shareholders at Annual General Meetings seem to concentrate more on statutory issues than on the goals of management. It is comforting, though, that directors must still approve key issues instrumental to the creation of economic wealth, such as the capital budget and financing decisions. The main methods employed by firms and their shareholders to overcome the agency problem are performance driven share and bonus schemes.

It is proposed that a performance measure such as Economic Value Added can and should be used to overcome the agency problem to benefit both shareholders and management.

**Key Words**

*Agency Problem*
*Economic Value Added*
*Performance incentive schemes*

# 1 The nature of the agency problem

## 1.1 Introduction and objective

Until approximately 1870, management and ownership of enterprises were vested in the same person, the capital provider (Lambrechts 1992:27). The emergence of large enterprises, especially the public company as a form of enterprise, was however characterised by a shift to separation between management and ownership of the enterprise. Owners appoint professional managers to manage their companies. It is this separation between ownership and management which forms the basis of the so-called agency theory. The shareholders' role become increasingly more passive while management has a reasonably free hand to pursue goals that may not necessarily correspond with those of the shareholders of the firm.

The goal of this article is to establish if (and to what extent) the agency problem exists among companies listed on the Johannesburg Stock Exchange. It also addresses various agency costs and investigates methods of overcoming the agency problem.

## 1.2 The agency problem

An agency relationship exists between the agent (management) and the principal (capital providers or owners) of the firm.

If both the agent and the principal are wealth maximisers (as we assume all rational people to be) then the possibility of conflict arises. The agent can and will take action to maximise his/her own wealth, and this action may not necessarily be in the best interest of the principal. If there is a difference between the goals management pursue and than of the owners (shareholder wealth maximisation), one can deduce that the agency problem is present. This aspect is investigated in the empirical section of the study.

Cohen and Uliana (1990:8) mention several examples of costs or actions by management which can give rise to excessive (more than normal) or unnecessary costs and which stem from this conflict situation:
- ☐ excessive levels of management remuneration;
- ☐ shirking (neglect of duty);

267

- [ ] the appropriation of corporate resources in the form of excessive levels of perks;
- [ ] avoiding investing corporate resources in potentially profitable ventures to the detriment of the shareholders;
- [ ] the pursuit of sales growth at the expense of profit or shareholder wealth;
- [ ] empire building by managers;
- [ ] employee welfare objectives; and
- [ ] manipulation of dividend policy at the expense of shareholder wealth creation.

If any of these costs are observed in an enterprise one can deduce that an agency problem is present.

Where the owner of an enterprise also attends to the management of the enterprise there can be no conflict between goals and therefore no agency problem exists. The more ownership is vested in people who are not directly involved in the management of the enterprise, the greater the possibility of conflict. It is therefore necessary to investigate the ownership structure of a company. This is done in the empirical part of the study.

## 1.3 Agency cost

To ensure that the goals of management correspond with those of shareholders, shareholders can institute certain incentive measures or monitoring steps. This does, however, have certain cost implications, with an accompanying detrimental effect on the wealth of the owners. The cost is highest if all management actions are monitored. This type of cost should only be incurred if the benefits to be derived are greater than the cost incurred.

Brigham and Gapenski (1993:21) define agency cost as all costs borne by shareholders to encourage managers to maximise shareholder wealth rather than act in their self interest.

Types of agency cost which can be identified include monitoring (for example auditing), structuring, opportunity and guarantee or insurance cost. In the empirical part of this study, certain aspects of these costs are further discussed together with the empirical results of their occurrence amongst the participants.

268

# 2 Methods to overcome the agency problem

## 2.1 Introduction

According to Brigham and Gapenski (1993:21), agency cost is low if the total remuneration of managers is linked to the market value of the company's share price. If the share price increases, both management's wealth and that of shareholders increase. There are, however, several factors beyond the control of management which influence the share price of the enterprise and which impair the affectivity of such a scheme. These factors are discussed briefly in Section 2.2 below.

Taking the above into consideration, the solution appears to be a shareholders' wealth-based incentive scheme (low agency cost) with some degree of monitoring (high agency cost).

## 2.2 Performance based incentive measures

Even before the classic work of Jensen and Meckling in 1976, on the implementation of the agency theory in the sphere of financial management, several studies were conducted regarding the connection between management's remuneration and the market price of a company's shares. Several measures are used to evaluate managers' performance. Some of the most common are sales, profit, current value of expected cash flows and value added.

According to Masson (1971:1286), linking a manager's remuneration to the share price, has two benefits. Firstly, managers should then act and make decisions in the best interest of the shareholders. Secondly, the stock exchange plays a reasonably effective role in the capitalisation of the future net income of the enterprise as represented by the share price. Managers should therefore concentrate on the net present value of the enterprise which is in turn closely link to shareholder's wealth.

### 2.2.1 Share option schemes

A share option scheme can be implemented by the enterprise in several ways and on several conditions. The basic principle is that an employee (usually at management level) has, as part of his/her remuneration, the right to purchase at a fixed price a number of the shares of the enterprise. This scheme is based on the principle that a participant will act and make decisions which will have a favourable influence on the share price of the enterprise. The reality is, however, that there are a number of factors which influence the share price of

an enterprise over which management has no control. Examples of these are interest rates, the general state of the economy (business cycles), foreign exchange markets, political activities, war and rumours.

Share option schemes, based purely on the share price as a yardstick, are thus regarded as a relatively poor incentive measure to encourage managers to act in the interest of the shareholders.

### 2.2.2 Performance-based share options

As a result of the shortcomings of "ordinary" share option schemes, there is an increasing move towards performance-related share options (Brigham and Gapenski (1993:21).

In terms of such schemes, managers in the enterprise are remunerated (with shares or cash) but on the basis of specific performance measures such as earnings per share, return on assets (ROA) or return on equity (ROE). The difference between the ordinary share option scheme and performance-based share option schemes, is that some objective yardstick is used in the latter instance in order to reward the managers.

It must be borne in mind that accounting-based measures can be subject to manipulation by management. Rappaport (1986:19) and Stewart (1990:35) reported in detail on the shortcomings of earnings per share and other accounting ratios in evaluating company performance. (It falls beyond the scope of this study to elaborate further on this topic, but the reader who is interested in furthering his knowledge on this subject can, in addition to the above-mentioned references, find an interesting discussion supporting the matter from Stern (1994:39)).

### 2.2.3 Value creation as performance measurement

One of the most important aspects in the evaluation of the performance of an enterprise is the creation or destruction of value to the capital employed. This value can also be used as a measurement to evaluate management. Economic Value Added (EVA) was developed and popularised by Joel Stern and Bennet Stewart (of the New York consultancy firm, Stern Stewart) over the past decade. Stewart formalised this philosophy and principals in his book, *"The Quest for Value"*.

EVA is a measure employed to evaluate what management has added to (or destroyed of) the total capital of the enterprise. Milunovich and Tsuei

(1996:106) stressed the importance thereof as a management tool in the following statement: "EVA instills capital discipline by forcing managers to consider the actual cost of the capital they employ. Thus, EVA encourages managers to act as owners."

The EVA of an enterprise is calculated as follows (Stewart 1990:137):

> EVA = Net income after tax
> **PLUS** certain adjustments (for example research and development costs)
> **LESS** [ weighted average cost of capital x capital ]

A remuneration scheme for managers can be based on the change in EVA (upwards or downwards) of an enterprise over a given period.

There are basically three ways whereby an enterprise's EVA can be improved (Stewart 1990:225):

1    increase net income without using more capital;
2    invest more capital as long as the return thereon is greater than the cost of that additional capital; and
3    liquidate capital or projects where the return is less than the cost thereof.

EVA is an all-embracing measure for the value which the managers of an enterprise add or deduct from the capital employed. To link managers' bonuses with changes in EVA links managers' remuneration directly to changes in shareholders' wealth (Stewart 1990:233). In this way, the management of a company is "forced" to set their goals and actions on economic shareholder wealth creation. The empirical section of this study reports on the number of companies that uses the EVA financial management system.

## 2.3 Shareholder control and interference

Cohen and Uliana (1990:9) defined "control" as the capability to elect the board of directors of a company. Even though control does not entail active decision-making in the enterprise, it does cover the taking of fundamental decisions such as the appointment of top management.

Shareholders can influence the company's management in two ways. Firstly, they can influence management directly as to how the company should be

271

managed. Secondly, any shareholder can make a proposal which is voted on at the annual general meeting (AGM) (Brigham & Gapenski 1993:23).

The most fundamental change which shareholders can effect is to change the board of directors. A board of directors which is controlled by management (too many executive directors as opposed to non-executive directors) is regarded as a weak link in the shareholder - management relationship.

The empirical section investigates shareholders' attendance at the AGM, the percentage ownership of directors as well as the matters on which both shareholders and directors vote at the AGM.

## 2.4 Threat of dismissal

In the past it seldom happened that a senior manager or chief executive officer was dismissed by shareholders. The reason for this was possibly that the ownership of a great number of companies was dispersed, as well as the fact that the agency problem was only brought to the attention of shareholders (and management) over the past two decades.

Nowadays, the reasons given for dismissal are shifting more and more away from "bad health" and "personal reasons" to "on the request of the directors" (Brigham & Gapenski 1993:25).

## 2.5 Threat of take-overs

Cohen and Uliana (1990:8) quoted an article by Jensen and Ruback that argued that the threat of a take-over serves to monitor the actions of management. If the actions or decisions of management decrease the future earnings or value of shareholders, the share price usually decreases as well. In some instances, the company can become a take-over target. If the management of such a company is replaced, the move can benefit the shareholders.
The threat of take-overs can thus serves as an external control mechanism which ensures that the decisions and actions of management maximise shareholders' wealth.

# 3. Empirical study

## 3.1 Research methodology

The primary aim of this study was to investigate the possible existence of the

agency problem at companies listed on the Johannesburg Stock Exchange. Furthermore, an investigation was conducted into the methods which are used to solve this problem.

Of the three most common methods used in the collection of data, namely personal interviews, telephonic interviews and postal surveys (Cooper & Emory 1995:287), the method of a questionnaire completed through postal survey and personal interviews was selected as the most suitable method for the study taking into account the costs involved as well as the availability of the participants. The population consisted of the listed companies on the Johannesburg Stock Exchange. When the survey was performed during September and October 1996, 577 companies were listed.

## 3.2 Questionnaire design

The questionnaire was in English. It consisted of only twelve questions on the front and back of an A4 page. The questions were phrased in such a way as to investigate both the possible existence of the agency problem and to deduce possible solutions to the problem.

The aim of the first two questions was to determine whether the agency problem does exist. The first question more specifically addressed the measurable goals which management pursue. The answers to this question indicate whether the goals management pursue was compatible with the firm's goal (wealth maximisation) and thus in the interest of the shareholders.

The next three questions attempted to ascertain to what extent shareholders try to control or monitor the decisions and actions of managers. The last seven questions investigated and analysed different aspects of agency cost.

## 3.3 Response

A total of 577 questionnaires were sent out. The following response was achieved:
☐   45 questionnaires were returned.
☐   In addition, personal interviews were conducted with 28 financial managers (part of the original population of 577) who had not taken part in the postal survey. This was done in order to increase the low response rate from the postal survey.
☐   7 enterprises indicated that it is company policy not to participate in any surveys.
☐   A total of 66 (45 + 28 - 7) (11%) responses could be used in the analysis.

273

It is important to note that some individual questions on the questionnaires returned, were not answered. Therefore the actual number of responses which could be used were in some instances less than 66.

## 3.4 Analysis of the results

The results are discussed **per question**, focusing on the aim, results and conclusion drawn.

### *Question 1*

Aim:
The aim of this question was to determine which goals the management of a company pursue. Goals which do not coincide with the maximisation of shareholders' wealth, are an indication of the existence of the agency problem.

Results:
The response of only 61 companies could be used for analysis. Value maximising is supported by 27 (44%) of the 61 companies while 36 (59%) indicated that improvement or growth in earnings is pursued as a goal. It was disappointing that only 4 companies specifically stated EVA as goal. Operating efficiency plays another major part (25%) while business ethics enjoys the lowest priority. Other goals include vision, cash flow, occupational safety and the increase in net assets.

Conclusion:
It appears from the response that a significant number of companies pursue goals which cannot be reconciled with an increase in shareholders' wealth. It can safely be assumed that the problem is aggravated by the apparent belief that an increase in earnings and return on equity (ROE) automatically result in an increase in shareholders' wealth. In the discussion on the theory it has been indicated that this might not necessarily be the case. It can therefore be concluded that an agency problem exists among a significant number of the respondents.

### *Question 2*

Aim:
The aim of this question was to determine how many shares are in the possession of management and employees and how many are owned by outsiders (who have nothing to do with the day to day running of the company).

The greater the percentage in the hands of outsiders, the greater the probability that the goals pursued by management do not necessarily correspond with the goals of the shareholders. From this it can be deduced that if the board of directors, as strategic executive organ of the company and representative of the shareholders, does not own a substantial percentage of the issued share capital, the possibility of conflict between the goals of the two groups increases.

Results:
Of the 66 companies' responses, 60 could be used for analysis. In 42 (70%) of the 60 companies, between 90% and 100% of the shares are owned by outsiders, while in 49 (82%) of the instances, the percentage shares owned by outsiders is between 70% and 100%. In two instances, the percentage of share ownership by outsiders in the company is between 50% and 100%.

An analysis of the percentage of shares owned by the board of directors, excluding the managing director, shows the following:

| | | |
|---|---|---|
| 0% | - | 10 companies |
| less than 1% | - | 23 companies |
| between 1% and 4% | - | 11 companies |
| between 4% and 30% | - | 6 companies |
| more than 30% | - | 10 companies |

Conclusion:
The relatively low percentage of shares owned by management and employees enhances the possibility for occurrence of the agency problem. This deduction is further supported by the low percentage of shares held by the directors in general.

*Question 3*

Aim:
The most direct way in which a shareholder who is not employed by the company nor serves on the board of directors can protect his\her interests, is by exercising his\her vote at the annual general meeting (AGM) which is compulsory in terms of the Companies Act. This question attempted to determine to what extent the shareholders employ this mechanism. If the voting percentage of shareholders at the AGM is low, this may also indicate a lack of knowledge among the shareholders about the existence of the agency problem.

Results:

The response of 65 companies could be used. The average percentage vote not exercised, amounts to 27%. The modus interval is 31% to 40% with 17 companies in this category. Four companies indicated that more than 90% of the votes are not exercised. One company indicated that 100% of the votes are not exercised!

Conclusion:

From these results it can be deduced that the average percentage of votes which are actually exercised is relatively high. Other factors which were not investigated in this study can naturally also play a role in the explanation of the voting percentage (high or low) at the AGM. For example, shareholders of companies which deliver good results year after year were most probably less inclined to use their vote or to attend the AGM at all.

## *Question 4*

Aim:

The aim of this question was to test the issue of control or monitoring over the actions and decisions of management by looking at the decisions which require the approval of the board of directors. The response to this question is to be read in conjunction with the response to Question 2, which investigated what percentage of the directors were also shareholders, which in turn can be an indication of the measure of control.

Results:

The response of 54 companies could be used. A total of 82 answers were obtained and these were divided into 8 categories. Of the 8 categories, it is mainly the capital budget, take-overs and mergers as well as the financing decision which must be approved by the board of directors. These three aspects were shown by 67%, 41% and 20% respectively of the 54 respondents to be decisions that need approval by the board of directors.

Conclusion:

In the literature study, it was shown that EVA is arguably the best method to determine shareholders' wealth. There are three main factors which determine the EVA of a company, namely adjusted net income after tax, capital employed and the cost of capital. It follows then that the aspects which influence these factors, namely the capital budget and the financing structure of the company, are the most important monitoring aspects according to the board of directors. It can thus be deduced that the directors can exercise a relatively important and efficient measure of control over shareholders' wealth.

## *Question 5*

Aim:
The aim of this question was to investigate further the aspect of control over the actions or decisions of management. In this instance, aspects for which approval from the shareholders is required were investigated.

The results of this question must be read in conjunction with the response to Question 3, which analysed the percentage of shareholders that exercised their vote (and thus their control action, even if it is small) at the annual general meeting.

Results:
The response of 59 companies could be used. 17% of the respondents indicated that none of the decisions had to be approved by the shareholders, while a further 17% indicated the same decisions which were supplied as answers in Question 4 (directors' approval). The balance, an overwhelming 66%, indicated that actions for which the approval of shareholders are required are stipulated by statute. This includes decisions such as share issues, remuneration of directors, appointment of auditors and directors, aspects regarding the Companies Act and the rules of the Johannesburg Stock Exchange.

Conclusion:
It appears that a very small percentage of the shareholders monitor or try to control those actions of management which can influence shareholders' wealth. Although the results of Question 3 indicated that the majority of the shareholders exercised their vote, it appeared to be mostly the statutory actions which enjoy attention. Strategic decisions are left to the directors, which emphasises further the separation between shareholders (ownership) and management. There is thus a considerable potential for the agency problem to arise.

## *Question 6*

Aim:
The aim of this question was to determine whether opportunity cost (see Section 1.3) occurs as a result of a complex, slow or cumbersome decision-making structure.

Results:
The response of 61 companies could be used. Seventeen (28%) of the respondents answered "yes" to this question. Unfortunately none of the companies responded to the request to supply examples.

277

Conclusion:
It is very difficult to determine or quantify opportunity cost. Because some of the respondents were of the opinion that opportunities in the market are lost as a result of the decision-making hierarchy, it can, however, be deduced that this cost does occur.

## *Question 7*

Aim:
The aim with this question was to determine to what extent incentive schemes, which attempt to harmonise the actions of management with the goal of shareholders' wealth, are used by companies.

Results:
The response of 65 companies could be used. It appeared clearly from the completed questionnaires that two main groups of incentive schemes are used, namely a share option scheme and a bonus scheme. The bonuses are normally based on some performance measure such as net income. Of the companies, 4 indicated that EVA is used as a measure.

Only 9 of the companies have no such scheme in place. Of the companies, 24 use a share option scheme, 13 use a bonus scheme while 19 use both.

Conclusion:
The occurrence of incentive schemes is very high and it thus appears to be one of the most popular methods by which the shareholders, as owners, try to circumvent the agency problem.

## *Question 8*

Aim:
The aim of Question 8 was to determine the occurrence of exceptional top management fringe benefits.

Results:
The response of 62 companies could be used. Only 5 of the companies answered affirmatively to this question. These benefits were comprised of, *inter alia*, increased provident fund contributions and increased pension fund contributions.

Conclusion:
Based solely on the completed questionnaires, it appears that this phenomenon does not manifest itself widely. The fact that so few of the respondents

278

answered in the affirmative to this question may be attributed to the fact that this type of expenditure is, on the one hand, not easily noticed or measurable and, on the other hand, that top management does not want to popularise (or acknowledge) the existence of such expenditure.

### Question 9

Aim:
Donations by companies to their favourite welfare organisations, although surely sociably acceptable, does not necessarily contribute to shareholder's wealth. It is very difficult to predict or quantify whether donations will have a positive or negative financial effect for a company or its shareholders. It falls outside the scope of this study to discuss this matter in detail. This question merely attempts to at least establish the extent of donations among the participant companies.

Results:
The response of 45 companies could be used. Of the companies, 32 (71%) indicated that donations amount to less than 1,5% of their net income, while 4 companies regarded their donations as insignificant. Of the companies, 9 designated a rand value of which only two amounts were substantial, namely R5 million and R15 million respectively.

Conclusion:
It does not appear as if donations play a significant role in the South African situation. The maximum allowable deduction for tax purposes (5% of net income), coupled with the fact that donations may only be made to approved institutions probably plays a direct role in what amounts are donated.

### Question 10 and 11

Aim:
The internal and external audit is the most important instrument for purpose of monitoring the actions of management see Section 1.3). Any (agency) cost should only be incurred in so far as the benefits exceed the cost thereof. This question attempts to test at least management's opinion in this regard.

Results:
Of the 61 companies which answered this question, 13 (21%) were of the opinion that the benefits of an internal audit do not justify the cost thereof.

Of the 61 respondents, 26 (43%) were of the opinion that the benefits of an external audit do not justify the cost thereof.

279

Conclusion:

Management is, in general, less critical regarding internal audits than they are in respect of external audits. The external audit is obviously compulsory, according to statute, which may, in the view of management, cause the nature and extent thereof to be less efficient. Although difficult (if not impossible) to quantify, it can be safely assumed that the (agency) cost of both the internal and external audits will be less than the potential benefits of this monitoring action undertaken on behalf of the shareholders. It is, however, significant that there is such a large measure of negative perceptions regarding the whole audit process, especially if one takes into consideration the fact that the questionnaires were mainly completed by financial managers. This is definitely an aspect which requires further research, especially from the viewpoint of the shareholder.

### *Question 12*

Aim:

Shareholders can protect themselves against a reduction in wealth as a result of the actions of management by letting the company take out fidelity guarantee insurance. This question attempts to determine the extent of this type of insurance.

Results:

Of the 51 companies that answered the question, 16 have no such insurance; 35 companies do have such insurance, but indicate the cost thereof to be insignificant or less than 1% of the net income; in 5 instances, the value amounted to more than 1%, while 4 companies indicated a rand value. The highest percentage was 10% and the highest rand value was R25 million.

Conclusion:

This type of insurance is fairly common and amounts, in certain instances, to substantial amounts.

## 4 Summary

In conclusion, it appears clearly from the responses to Questions 1 and 2 that the agency problem does exist in a substantial number of cases. The goals of management differ from the goals of shareholders. This deduction is further supported by the fact that a relatively small percentage of the issued share capital is owned by management and the board of directors.

At first glance it appears that shareholders do try to address the agency problem by the relatively high voting percentage at the annual general meeting,

280

but after further investigation it appeared to be mostly statutory aspects which enjoyed attention at the annual general meeting. What is, however, comforting is the fact that the board of directors must approve certain aspects which are instrumental in the creation of economic wealth, such as capital budgets and the financing decision.

The most important method used to overcome the agency problem is performance- based shares, bonus schemes and external audits.

It was furthermore found that agency cost appears in the form of opportunity cost, donations and insurance. Almost no exceptional top management benefits appeared.

EVA as goal and a performance measure is found in relatively few companies in South Africa. Some of the world's biggest conglomerates (Coca Cola, Quaker Oats and Briggs and Stratton) use the EVA financial management system with its many advantages, as indicated in the discussion on the theory. It falls beyond the scope of this study to entertain a detailed discussion on EVA, but EVA can provide the link between management and the shareholders so that agency cost is not only minimised, but the agency problem in itself is also greatly reduced.

# Bibliography

Brigham, E.F. & Gapenski, L.C. 1993. _Intermediate Financial Management_, fourth edition, The Dryden Press, Orlando.

Cohen, T. & Uliana, E. 1990. _An evaluation of corporate ownership structures on employee, management and shareholder compensation for JSE companies_, De Ratione, Vol.4(1), pp.7-14.

Cooper, D.R. & Emory, C.W. 1995. _Business Research Methods_, fifth edition, Richard D Irwin.

Jensen, M.C. & Meckling, W.H. 1976. _Theory of the firm: managerial behaviour, agency costs and ownership structure_, Journal of Financial Economics, Vol.3, pp.305-360.

Lambrechts, I.J.(Ed) 1992. _Financial Management_, fourth edition, JL van Schaik, Pretoria.

281

Masson, R.T. 1971. *Executive motivations, earnings, and consequent equity performance*, Journal of Political Economy, Vol.79(2), pp.1278-1292.

Milunovich, S. & Tsuei, A. 1996. *EVA in the computer industry*, Journal of Applied Corporate Finance, Vol.9(1), pp.104-115.

Rappaport, A. 1986. *Creating Shareholder Value*, The Free Press, New York.

Stern, J.M. 1974. *Earnings per share don't count*, Financial Analysts Journal, July-August 1974, pp.39-45.

Stewart, G.B. 1990. *The Quest for Value. A Guide for Senior Managers*, Harper Collins Publishers.

*JH Hall*
*Senior Lecturer*
*Department of Business Management*
*University of Pretoria*
*Pretoria*

# Comparing Varieties of Agency Theory in Economics, Political Science, and Sociology: An Illustration from State Policy Implementation*

Edgar Kiser

*University of Washington*

*As rational choice theory has moved from economics into political science and sociology, it has been dramatically transformed. The intellectual diffusion of agency theory illustrates this process. Agency theory is a general model of social relations involving the delegation of authority, and generally resulting in problems of control, which has been applied to a broad range of substantive contexts. This paper analyzes applications of agency theory to state policy implementation in economics, political science, and sociology. After documenting variations in the theory across disciplinary contexts, the strengths and weaknesses of these different varieties of agency theory are assessed. Sociological versions of agency theory, incorporating both broader microfoundations and richer models of social structure, are in many respects the most promising. This type of agency theory illustrates the potential of an emerging sociological version of rational choice theory.*

Rational choice is not just a microlevel theory, but also consists of a set of general models of social structure and social relations (Becker 1976; Schelling 1978; Friedman and Hechter 1988; Hechter and Kanazawa 1997; Kiser and Hechter 1998). Agency theory is one such general model,[1] with roots in both the classical work of Max Weber and the "new institutionalism" in economics that is now also being used in political science and sociology. An agency relation is one in which a "principal" delegates authority to an "agent" to perform some service for the principal. Agency relations exist in a wide variety of social contexts involving the delegation of authority, including clients and various service providers (doctors, lawyers, insurance agents), citizens and politicians, political party members and party leaders, rulers and state officials, employers and employees, and stockholders and managers of corporations (in all cases, the former is the principal and the latter the agent). The key feature of all agency relations is that once principals delegate authority to agents, they often have problems controlling them, because (1) agents' interests often differ from theirs, and (2) agents often have better information about their actions than do principals. Agency theory focuses on the ways principals try to mitigate this control problem by selecting certain types of agents and forms of monitoring their actions, and by using various amounts and types of positive and negative sanctions.

This paper critically analyzes the intellectual history of agency theory in three disciplines—economics, political science, and sociology. The substantive focus will be on

[1]Other rational choice models include optimal location theory (Downs 1967), group solidarity theory (Hechter 1987), and the family of models usually called "game theory."

state policy implementation (the relationship between rulers/politicians as principals and state officials/bureaucrats as agents) because scholars in all three disciplines have applied agency theory to this topic.[2] The most general issue in this literature is how do the leaders or rulers of states control the many state officials to whom they delegate authority to implement policy.

The substantive theoretical goals of the paper are: (1) to compare the strengths and weaknesses of the different forms of agency theory that have developed in the three disciplines, and (2) to evaluate the overall success of agency theory as a model of state policy implementation.[3] In addition, some more metatheoretical questions will be addressed, including: (1) what can the spread of agency theory tell us about the process of intellectual diffusion—is it best described as a case of "economic imperialism" (in which economic ideas are seen as invading and replacing those of other disciplines)?; (2) what can we learn about standard disciplinary modes of explanation by analyzing the same theory in different disciplinary contexts?; and (3) what are the implications for the future of rational choice theory in disciplines like political science and sociology?

I begin by discussing the development of formal agency models in economics that emerged as part of the "new institutionalism" (Ross 1973; Jensen and Meckling 1976). Although these models have been applied mainly to economic organizations, there have been a few recent applications to state policy implementation (Rose-Ackerman 1978; Klitgaard 1988; Demski and Sappington 1987; Johnson and Libecap 1989). This offers an opportunity to study the use of agency theory in a new substantive domain within its disciplinary "home turf." Current work on state policy implementation using agency theory in political science (Weingast and Moran 1983; Bendor and Moe 1985; Wood 1988; Kiewiet and McCubbins 1991) has drawn mostly on economic models, although they have been transformed in several ways by the disciplinary context of political science. I then turn to the classical sociological origins of agency theory[4] in the work of Max Weber ([1922]1968).[5] Work using agency theory in sociology (Hamilton and Biggart 1980, 1984; Kiser and Tong 1992; Gorski 1993; Kiser and Schneider 1994; Kiser 1994; Adams 1996) is then analyzed as a synthesis of economic and Weberian ideas, producing a sociological version of agency theory that in different ways tries to combine the breadth of Weber with the precision of economic models. The way agency theory has been used in political science and sociology is much different from its use in economics, indicating that the "economic imperialism" metaphor is not a good description of the intellectual diffusion process in this case. The evolution of agency theory as it has diffused across disciplines suggests that rational choice theory is becoming increasingly separated from the discipline of economics and is dramatically transforming as a result.

As a rough guide to the detailed discussions of different forms of agency theory, Table 1 provides a short summary of some of the main characteristics of the theory in economics, political science, and sociology. Since agency models in all disciplines deal with monitor-

[2] Agency theory has primarily been applied to the analysis of economic organizations. However, since almost all of this work has been in economics, this substantive focus would not allow for comparisons across disciplines.

[3] A systematic evaluation of agency theory as a model of state policy implementation would require comparing it to other theories, a task beyond the scope of this paper. Here we explore the extent to which agency theory has been able to adequately respond to recent criticisms (see especially Perrow [1990]).

[4] The other main classical sociologist who used a proto-agency theory framework was Roberto Michels ([1915]1959). His analysis focused on the relationship between political party leaders (as agents) and party members (as principals). He argued that as parties grow in size, informational asymmetries between party members and leaders increase, resulting in a decreasing ability for members to control their leaders (his famous "iron law of oligarchy"). I focus on Weber mainly because he is more interested in the issue of state policy implementation.

[5] In spite of his innovative theoretical analysis of agency problems in political institutions and the unprecedented range of his historical examples, Weber's work had practically no influence on the development of agency theory in economics and very little influence in political science. The rediscovery of Weberian ideas and themes is one of the main things that makes sociological analyses using agency theory unique.

Table 1. Main Characteristics of Agency Theory in Economics, Political Science, and Sociology

|  | Economics | Political Science | Sociology |
|---|---|---|---|
| Microfoundations | Parsimonious, rationality and self-interest;[a] also include risk | Parsimonious, rationality and self-interest[b] | Ranges from parsimonious to broad; often include values |
| Meso Level: Organizational Structure | Sparse or absent;[c] some include multiple agents | Sparse; some include multiple principals | Several ideal types of organizational structure; include both multiple principals and multiple agents |
| Macro Level: Structural Context | Sparse or absent | Sparse or absent | Very full; all incorporate "material" structure, many include culture as well |

[a]Rose-Ackerman (1978) and Klitgaard (1988) are exceptions, since they use broader microfoundations.
[b]Banfield (1975) is one of the few political scientists to use broader microfoundations with agency theory.
[c]Rose-Ackerman (1978) is an exception—she uses models of organizational structure.

ing, sanctioning, and (usually) recruiting of agents, the table focuses on the features that make each most distinctive: their microfoundational assumptions, and the way they usually model the organizational and structural contexts in which agency relations are embedded.

## AGENCY THEORY IN CONTEMPORARY ECONOMICS

Adam Smith ([1776]1976:700) was well aware of agency problems in organizations, including those arising from the separation of ownership and control:

> The directors of such companies, however, being the managers rather of other people's money than of their own, it cannot well be expected, that they should watch over it with the same anxious vigilance with which the partners in a private copartnery frequently watch over their own. . . . Negligence and profusion, therefore, must always prevail, more or less, in the management of the affairs of such a company.

However, economists did not elaborate Smith's insights on this topic for almost two centuries. The firm remained a "black box" in economic models; until very recently economic organizations were reduced to "production functions." As Jensen (1983:321) describes the situation before the rediscovery of agency problems, "[i]n most economic analyses, the firm is modeled as an entrepreneur who maximizes profits in an environment in which all contracts are perfectly and costlessly enforced." Stiglitz (1987:242) notes that this left neoclassical economics with no theory of the relationship between incentives and performance within firms: "The standard theory was based on the assumption that what action

285

the 'principal' willed his agent to perform was perfectly known, and that the action could be perfectly and costlessly monitored. Neither assumption is plausible."

Agency theory emerged as a result of two dramatic changes in economics. First, several economists began to redefine their discipline, from one focused on the object of study (economic relations) to one based on an economic theory (or "approach") that could be applied to any substantive domain. As a consequence, many economists branched out into new topics, including politics (Arrow 1951; Becker 1957, 1976; Buchanan and Tullock 1962; Downs 1957), and "public choice" theory was born. Second, the "new institutionalism" (sometimes called the "new economics of organization") in economics began to open up the black box of the firm, with a new focus on property rights (Cheung 1969; North 1981; Barzel 1989), transaction costs (North and Thomas 1973; Williamson 1976, 1985), and agency problems (Ross 1973; Jensen and Meckling 1976). When public choice and agency theory came together, one of the products was the application of agency theory to the study of state policy implementation (Rose-Ackerman 1978; Klitgaard 1988).

There are several good summaries of economic theories of agency available (MacDonald 1984; Eisenhardt 1989; Petersen 1993), so my comments here will be brief and focused on the general process of theoretical development. In the most general terms, economic agency theorists have used parsimonious microfoundations (rationality, self-interest, wealth maximization) and simple models of organizations (theoretically viewed as a nexus of contractual relations, but in practice usually only one principal and one or two agents) to explore the consequences of variations in monitoring capacity and incentive schemes for the performance of agents and thus for the overall efficiency of organizations.

Agency theory in economics was initiated by the work of Berhold (1971), Ross (1973), and especially Jensen and Meckling (1976) as a way to address problems of control that arise as a result of information asymmetries between agents delegated to carry out tasks that affect the welfare of the principals who delegated authority to them.[6] Berhold (1971) focuses on the way risk preferences affect contract choice (the more risk averse agents are, the more they want fixed salaries), and the effects of different monetary incentive schemes on performance. These two issues have continued to be the main foci of economic theories of agency. Jensen and Meckling (1976) develop the model much more fully, and apply it to the agency problem discussed by Adam Smith, the separation of ownership and control.

The role of risk is an important innovation in the economic theory of agency. Economists recognized that agency relations involve not only problems of control but issues of risk-sharing as well. If all actors are assumed to be risk neutral, then residual claimancy (who gets the "residual" profit or loss from a joint activity, as opposed to getting fixed share regardless of the outcome) is allocated between principals and agents on efficiency grounds. However, risk adverse actors will often reject a contract in which they bear most of the risk even if it is efficient. Risk adverse agents essentially purchase insurance by accepting fixed salary contracts with lower total value but less variation. Therefore, taking risk into account allows economists to explain the existence of "inefficient" contracts (more strictly speaking, contracts that would be inefficient if all actors were risk neutral). However, the advantage of including risk in the theory is often lost when the mathematical models are constructed—tractability issues often require assumptions of risk neutral agents. Even when risk preferences are included in the theory, they are never directly measured.

The selection of agents is important in economics arguments, since one of the things stressed by several analyses is that there are different "types" of agents, meaning that agents differ in their general levels of ability, effort, and honesty (this is often summarized as "high productivity" vs. "low productivity" agents). These differences are better known

---

[6]Some of the key issues about the nature of firms were raised earlier by Coase (1937), and some of the same topics discussed in agency theory were addressed in Arrow's (1964) work on "moral hazard."

to agents than to principals, another aspect of the information asymmetry between the two, which creates a problem for principals attempting to select agents: How can they tell the good ones from the bad? In most economic models, agent type is exogenous—the usual way of talking about this is that "nature draws a type for the agent" (Kofman and Lawarree 1996:122). However, economic agency theorists have made two main arguments about agent type. The first is that principals can affect the type of agents they get by the form of the contract they offer, since agents, knowing their own "type," will select the contract most beneficial to them. For example, low productivity agents will choose fixed salaries, whereas high productivity agents will choose piece-rate contracts. The second main argument comes from the related literature on "signaling" (Spence 1973). Sellers of high quality products (such as high productivity workers seeking employment) will attempt to reveal their "type" by undertaking some activity that is less costly to them than it would be to someone selling a low quality product. Long and inclusive warranties are an example in product markets; the willingness to pay a bond would be an example in the labor market for agents.

### Strengths and Limitations of Economic Agency Theory

The main criticism of economic agency theory by sociologists is the familiar general criticism of rational choice theory, that it is too sparse in its depictions of both the micro and macro levels. For example, Perrow (1990) rejects agency theory primarily because it does not model the context of behavior in a realistic manner and relies on an assumption of self-interest that he views as too narrow.[7] Petersen (1993:288–89), MacDonald (1984:429–30), and Eisenhardt (1989:71–72) all agree that agency theory as used by economists is too parsimonious, and that it can and should be broadened to include additional elements.

   One of the main determinants of the limitations of economic agency theory is the narrowness of its empirical scope. In sharp contrast to Weber's broad discussion of agency problems in many types of organizations in varied historical contexts (see below), agency theory in economics has focused on two main types of agency relations in economic organizations—between stockholders and managers, and between employers and employees. The empirical focus has usually been limited to economic organizations in the contemporary United States. Due both to this constrained empirical scope and the desire to create tractable mathematical models, economists have focused on a limited range of monitoring and sanctioning strategies. For example, they have reduced the issue of selection of agents to the problem of "adverse selection"—the tendency of individuals with substandard qualifications to self-select into positions as agents if principals use ineffective selection criteria or offer lower-than-average compensation. Economic agency theorists have also focused on a fairly limited range of incentives, usually various forms of salary payment (either fixed or dependent on output), and sometimes including the threat of dismissal as well. Since they have tended to explore agency problems within a narrow range of structural and organizational contexts, they have implicitly ignored a wide range of structural determinants of agency outcomes by holding structure constant.

   Another serious limitation of economic agency theory, following mostly from the severe parsimony of its mathematical models, is that it is an organizational theory without organizations. In fact, the typical "organization" in agency models until quite recently con-

---

[7]Perrow thinks that agency theory is politically biased since it ignores the possibility that principals (often employers) may act opportunistically to exploit agents (often employees)—in other words, it looks only at the sins of the powerless, not the powerful. This is part of a more general common criticism of rational choice theory by sociologists, that it contains an inherent conservative political bias. His claim about bias might be true for studies of firms, but it is clearly false for agency theory more generally. In many cases, agents are more powerful than principals (rulers relative to officials or citizens, doctors relative to patients). In these situations, agency theory focuses on how those with less power can control those with more.

sisted of only two actors, one principal and one agent.[8] The problem with using such simple models is that they are unrealistic even within the context of that narrow range of organizations of interest to most economists. Stockholders of corporations are multiple principals, and employees are multiple agents—and these facts are often critical for understanding agency relations (multiple employees can engage in collusive corruption, for example, or they can either report or hide the shirking of other employees). Recently, economists have begun to address the issues raised by multiple agents (Tirole 1986, 1992; Bac 1996; Kofman and Lawarree 1996)—for example, Tirole (1992:158) makes a compelling argument that centralization decreases collusion, since collusion depends on agent discretion, and that is reduced in centralized organizations—however, they have not yet expanded their models to include multiple principals[9] (in this respect, political science and sociological applications are more advanced).[10] This makes it difficult to apply economic models of policy implementation to most contemporary democratic states, since they are characterized by multiple principals.

One of the consequences of the recent expansion of economic models to include multiple agents is that they have rediscovered issues central to classical sociology. For example, Bac suggests that "[b]y committing himself to a probability of rotating the agents among various hierarchical positions, the principal or the designer of an organization can affect the discount factors of the members, hence indirectly keeping internal corruption under control" (1996:290). Although this is offered as a novel insight, it was discussed by Weber ([1922]1968:1043) since it was a tactic often used by medieval and early modern rulers. However, this does not mean that all economists are simply recapitulating Weber with no value added. Kofman and Lawarree (1996) analyze a three-level hierarchy with a principal, a supervisor, and an agent, and explore the possibility of collusion between the supervisor and the agent. They argue that collusion can be reduced by introducing redundancy—using a second supervisor to share monitoring tasks with the first. This is also a rediscovery of Weber's ([1922]1968:222, 995–97) analysis of collegial monitoring institutions, which he argues is a patrimonial form used to reduce corruption when monitoring problems are severe. However, Kofman and Lawarree go beyond Weber's analysis by better specifying the type of relationship that must exist between the two supervisors for collusion to be adequately deterred (a prisoner's dilemma situation). This is an interesting case of formal modeling being able to add to Weber's informal analysis.

*Economic Agency Theory Applied to State Policy Implementation*

One thing that all of the early economic agency theorists have in common is a tendency to highlight the generality of their models. For example, all of them note that agency theory could be used to study state policy implementation (Ross 1973:134; Jensen and Meckling 1976:309; Jensen 1983:321). Jensen and Meckling (1976:309) go on to say that "[t]he development of theories to explain the form which agency costs take in each of these situations (where contractual relations differ significantly) and how and why they are born

---

[8]Bac (1996:277) notes that "The literature examining the problem of motivating agents to be honest mostly takes the institutional context as given."

[9]They also tend to use fairly restrictive assumptions about how principals will act. For example, a common assumption in most economic agency models is that the principal can commit to following a particular monitoring and sanctioning strategy ex ante, and will not change it ex post (called the "Stackelberg assumption"). This may be realistic in some settings, but it is a poor approximation of the actual behavior of absolute monarchs relative to their administrative staffs, for example.

[10]Even when more than one agent is introduced, it is often in a superficial way. For example, Kofman and Lawarree (1996) note that most agency models that have introduced third parties as supervisors have assumed that the supervisor is simply a "third arm" of the principal with no independent or conflicting interests. They improve dramatically on these models by assuming that the supervisor is self-interested. However, they make other fairly drastic simplifying assumptions, such as that the supervisor can always perfectly monitor the agent.

will lead to a rich theory of organizations which is now lacking in economics and the social sciences generally." This tendency to stress the generality of theories is part of the culture of economics as a discipline, and no doubt it is one of the main reasons for the rapid diffusion of agency models from economics to other disciplines, and for their application within economics to state policy implementation.

Niskanen's *Bureaucracy and Representative Government* (1971) is probably still the most widely read and cited economic analysis of agency problems in state policy implementation.[11] Niskanen does not use formal agency theory, but he clearly frames the difficulties politicians have controlling bureaucracies as agency problems. He argues that one of the main threats to contemporary democracy is that elected politicians are losing power relative to appointed bureaucrats.[12] Niskanen clearly specifies the preferences of bureaucrats (agents in his model), which allows him to predict the direction of the deviation from politicians' (principals') preferences. He assumes that bureaucrats want to maximize the budget of the agency they control. They are often able to do this, in spite of the fact that it is contrary to the interests of principals, because they have better information than do politicians or voters about what budget level is necessary for their bureaucratic agency to adequately carry out its mission. As a consequence, the state grows larger and more expensive than either politicians or voters want it to be. Niskanen goes on to analyze various types of monitoring and sanctioning strategies that might mitigate the control problem. He even gets back to some classic Weberian solutions such as rotation—suggesting that rotating members on congressional committees might be a good way to minimize collusion between committee members and the agencies they oversee.

Although few economists have followed Niskanen's lead in studying agency problems in state policy implementation (the main exceptions are Rose-Ackerman 1978; Klitgaard 1988; Chandler and Wilde 1992),[13] some of the earliest applications retained some of the breadth of Weber and Niskanen while becoming more formal and analytical. Rose-Ackerman's (1978) fascinating study of corruption is perhaps the best example.[14] Although she does not draw explicitly on Weber's work, Rose-Ackerman is in several respects as close to Weber as she is to contemporary economists. She begins with the realization that models used to analyze economic organizations cannot be mechanically extended to the state: "since both politician and bureaucrat operate in distinctive institutional frameworks different from those of competitive theory, a simplistic application of market analysis is not sufficient" (1978:3). This implies that economic theory must be "tempered by a concern for the structure of government institutions" (1978:3). Therefore, unlike most (if not all) prior applications of agency theory, she differentiates among four types of organizational structure (fragmented, sequential, hierarchical, and disorganized), and examines the varying dynamics of corruption within each (ibid.:168).

In addition to taking organizational structure into account, Rose-Ackerman also looks at the importance of third parties and multiple agents, and even expands the usual rational choice microfoundations. Third parties are important in many agency relations, but they are usually ignored in economic agency models. For example, state policy implementation always involves third parties (taxpayers, regulated industries) and the actions of these third parties often effects implementation outcomes. Rose-Ackerman (1978:87) argues that the extent to which the beneficiaries of state policies are organized will affect the

[11] Downs (1967) also developed a pre-agency model of bureaucracy based on public choice theory, but unlike his *Economic Theory of Democracy* (1957), his work on bureaucracy has not been very influential.

[12] These Weberian themes (see below) are no accident—Niskanen has read Weber and cites many of his arguments. It is unfortunate that he is one of the last economists to do so.

[13] A recent summary of contemporary work on agency theory in economics (Kotowitz 1987:212) concludes that "[t]he consequences of moral hazard in political processes have largely been neglected by economists. . . . The theoretical tools of agency, contract, and game theory have yet to be fully employed in this area."

[14] Rose-Ackerman was the first economist to explicitly use agency theory to analyze state policy implementation.

success of policy implementation. The possibility of collusion between multiple agents is also addressed (ibid.:216–17). Perhaps most interesting for an economist, Rose-Ackerman is quite willing to consider the possibility of action not based simply on self-interested wealth maximization. She wants to "use economic analysis itself to refine the understanding of the place of morality in social life" (1978:229). "Economics does have a central role . . . in clarifying those situations where extraordinary appeals to personal integrity are not necessary for effective administration" (ibid.:218). In short, Rose-Ackerman seems to follow Weber in viewing the microlevel assumptions of rational choice theory as an ideal type useful not only for discovering regularities, but for highlighting deviations from rationality, as well (see also Downs [1967:88–89] on different "types" of agents).

Klitgaard's (1988) analysis of corruption in political institutions is one of the few to follow up on Rose-Ackerman's important effort to combine the rigor of mathematical economic models with some of the breadth of Weber. He too uses agency theory as his core model, and most of his analysis focuses on central issues of monitoring and sanctions, but he also adds many elements usually ignored by economists. Third parties are also central to his analysis; he notes that taxpayers can often be used to control corruption if they are encouraged to report these practices to higher authorities (1988:86). Since many forms of corruption hurt taxpayers, they have strong incentives to serve this monitoring function. Klitgaard also discusses the problems with collusion that arise when principals have multiple agents (ibid.:16–17). Although like Rose-Ackerman he primarily uses rational choice microfoundations, he tries to take cultural determinants of corruption into account, as well.[15]

Perhaps the most unfortunate aspect of the short history of agency-based analyses of policy implementation in economics is that works like those of Rose-Ackerman and Klitgaard have been so rare. The largest body of work in this area (still quite small) has focused on tax administration, and it has been much closer to mainstream economics in its theory and methodology.[16] Economists' models of tax administration began just as their models of firms had—they explored the relationship between the state and the taxpayer, leaving the internal dynamics of the state a "black box" (Reinganum and Wilde 1985; Scotchmer 1986; Border and Sobel 1987; Mookherjee and P'ng 1989; Graetz, Reinganum, and Wilde 1986).[17] As agency theory belatedly began to make some inroads in this area, several scholars started to disaggregate the state. Sanchez and Sobel (1993) and Cremer, Marchand, and Pestieau (1990) looked at the ability of politicians to control the IRS, and Demski and Sappington (1987) have done the same with the relationship between congress and regulatory bureaucracies. The limitations of this work are the general problems with economic agency theory discussed above—the empirical scope is very narrow and unrealistic assumptions abound due to requirements of mathematical modeling.[18] Although most of the formal work on tax administration does illustrate these classic shortcomings of the discipline, the studies done by Niskanen, Rose-Ackerman, and Klitgaard indicate that

[15] Another interesting paper in this tradition is by Johnson and Libecap (1989). They also note how the contemporary state is different institutionally from the firm—negative sanctions are weak and hard to enforce, and positive sanctions are also constrained by rigid salary scales. How then, can compliance be maintained in this agency relationship? They argue that agents are compelled to act in the interests of principals by bureaucratic civil service rules. Although these rules may in some sense be inefficient, since they add countless documentation requirements for all actions, they in fact enhance efficiency by lowering monitoring costs and making corruption more difficult.

[16] There has also been some work on the relationship between salary levels and performance in the public sector that is closely related to agency theory. Goldstein and Ehrenberg (1976) and Ehrenberg, Chaykowski, and Ehrenberg (1988) studied local political officials and found positive associations between pay and performance. More recently, Bates and Santerre (1993) have questioned the causal direction of this relationship—they find positive effects of performance on salary, but no effect of salary on performance.

[17] The same is true of the closely related literature on regulation (Stigler 1971; Peltzman 1976).

[18] Sanchez and Sobel (1993:367) assume away bribery, for example, and their only justification for this and other restrictive assumptions is that "[o]n a technical level, we have placed several ad hoc restrictions on the model." Cremer, Marchand, and Pestieau (1990:70) also admit to making many unrealistic assumptions "for the sake of tractability."

economic analyses can be much broader. It will be interesting to see if future economic applications of agency theory to policy implementation continue to develop in both of these directions.

## AGENCY THEORY IN CONTEMPORARY POLITICAL SCIENCE: CAN POLITICIANS CONTROL BUREAUCRATIC AGENCIES?

Contemporary political scientists have begun to use agency theory to explore one of the most significant issues Weber ([1922]1968) raised about contemporary politics: can elected politicians control the appointed bureaucrats who implement state policies, and if not, does the power of unelected bureaucrats threaten democracy? However, they have borrowed both theory and method from economic agency theory, not from Weber. One result of this is that their incorporation of structural embeddedness and heterogeneous microfoundations has been only partial.

The question of how politicians are able to control appointed bureaucrats has a long history in political science, including work on "iron triangles" (Truman 1951; Freeman 1955) and "regulatory capture" (Huntington 1952; Bernstein 1955).[19] The first person to systematically use agency theory to address this question was Banfield (1975). Like Rose-Ackerman, Banfield focuses on the ways in which political organizations differ from economic ones, and how this affects the agency relations within them. His use of agency theory is both broad and informal. His microfoundations are rational choice, but he also recognizes the importance of various deviations from standard rational choice assumptions, discussing both the possible "psychic costs" of corruption, and the role of preferences for power, glory, and "serving a good cause" (1975:588, 596). Banfield is also one of the first scholars to seriously discuss the problems created by multiple principals (ibid.:595). The existence of collective action problems among principals (in this case, usually congress or congressional committees) often makes monitoring more difficult and corruption in bureaucratic agencies higher. This central point is ignored in the economic literature, but plays a key role in political science. Banfield's role in political science is in many respects similar to Rose-Ackerman's in economics—both used agency theory broadly and informally[20] and raised several fundamental issues that others would later address using more formal models.

Although rational choice work in political science has focused most of its efforts on studying legislatures (Moe 1990:214, 222), there is now a fairly large literature applying agency theory to policy implementation. The focus has been on the general issue of "bureaucratic drift"—the tendency for the actions of bureaucratic agencies to "drift away" from the goals of the politicians trying to control them. The form of many of the arguments is partly functionalist: existing structures (procedures, monitoring systems, etc.) are explained by the agency problems they presumably mitigate (ibid.:224; Fiorina 1990:256).

Since the central problem in all agency relationships is information asymmetry, agency analyses naturally tend to focus on monitoring. One interesting finding in the political science literature is that there seems to be little direct monitoring of bureaucratic agencies by their political principals (Weingast and Moran 1983; Hammond and Knott 1996).[21] This discovery led some to conclude that Weber and Niskanen must be right, that bureau-

---

[19] "Iron triangles" refer to mutually profitable collusive relationships among interest groups, congressional committees, and the bureaucratic agencies these committees are supposed to control. "Regulatory capture" indicates a situation in which a regulated industry in fact controls (captures) the bureaucratic agency that is supposed to regulate it, in spite of attempts to limit this by congress.

[20] Some of Rose-Ackerman's (1978) analysis does use formal mathematical models, but most of it is informal verbal theory.

[21] There is of course some direct monitoring of bureaucratic agencies by organizations, including the Congressional Budget Office and the General Accounting Office (using hearings, investigations, and budget reviews), but scholars seem to agree that there is far too little monitoring to produce compliance in the absence of other factors.

cracies are indeed beyond the control of politicians—but this still leaves a critical question unanswered: why are politicians hardly even trying to monitor them? The answer many political scientists have converged on, especially in earlier work in the "congressional dominance" tradition (Weingast and Moran 1983; Weingast 1984), was that congress was in fact able to adequately control bureaucratic agencies (and thus, although this was never noted, that Weber's fear was unwarranted), but that they used means other than direct monitoring. As McCubbins and Schwartz (1984:166) note, direct monitoring is expensive, thus principals have strong incentives to find less costly strategies. One way to compensate for poor monitoring is to use stronger sanctions (Becker and Stigler 1974). Following this line of thought, several scholars (Weingast 1984; Weingast and Moran 1983; McCubbins and Schwartz 1984) have argued that the use of strong ex post sanctions (adjusting budgets using appropriations and reauthorizations bills) are key components of congressional control over bureaucracy. Kiewiet and McCubbins (1991:18) conclude that most of the causal arguments in political science agency analyses focus on the use of sanctions.

Even with strong sanctions, some monitoring is still necessary—how does congress know whose budgets to cut? Since most types of bureaucratic "drift" harm some interest groups or other citizens, these third parties have strong incentives to monitor bureaucracies and report problems to politicians. This type of reactive, third-party "fire alarm" oversight is much cheaper for politicians (the costs are paid by the third parties) than direct monitoring, and probably more effective (Weingast 1983; Kiewiet and McCubbins 1991:27–34). However, Weingast and Moran (1983:767) also note that third parties can hinder monitoring; for example, interest groups may collude with agencies to serve their mutual interests, which my include hiding some agency actions from politicians.

Administrative procedures can also be used to mitigate monitoring problems.[22] In a manner similar to Johnson and Lebicap (1989) (although there are no cross citations), several political scientists (McCubbins 1985; McCubbins, Noll, and Weingast 1987:254–55) have argued that various types of administrative procedures can mitigate informational disadvantages with reporting requirements ("red tape" thus serves a monitoring function), and can facilitate third-party monitoring by giving particular constituencies access to the agency. It might also be useful to think of additional administrative procedures as a negative sanction from the agents' perspective. Almost all administrative procedures add to the workload of agents without improving their ability to complete their tasks. Therefore, agents will be interested in avoiding the imposition of additional procedures. They may even comply with current rules in order to avoid having new ones added. Thus, administrative procedures could be used by congress as potential sanctions.

How should the preferences of politicians and bureaucrats be specified? Politicians are clearly the easier of the two, since there are abundant written records of their votes.[23] Thus Weingast and Moran (1983) and Moe (1985) are able to construct "ideological indexes" of the congress as a whole, and even for relevant committees and subcommittees. This type of detailed preference specification is very rare in rational choice models—the data availability in this context makes this a great research site for testing agency models. The preferences of bureaucratic agents are much more difficult to specify. Most contemporary political scientists reject Niskanen's (1971) assumption that the main goal of bureaucrats is to maximize their budget (or "slack"). Most of them follow economic models and focus on the effort made by the agent (Kiewiet and McCubbins 1991:33). Wood argues that "principal-agent models assume bureaucrats are passive, lazy, and calculative only to the extent that they want to avoid work" (1988:791). Adopting these microfoundations from economic

---

[22] Bawn (1995) notes that there is often a trade-off between technical efficiency and political control. Administrative procedures can enhance control, but they often decrease efficiency (red tape does have costs).

[23] There are problems with using voting records as indicators of preferences, however. They could reflect party discipline, or they could be strategic.

models has been problematic—this is an area in which future research could improve upon existing models.[24]

One of the major problems in the political science literature is that there is no clear consensus about who is the principal that is supposed to control bureaucratic agencies—is it the president, congress, an "enacting coalition" in congress, congressional committees, or some combination? All are aware that there are multiple principals, but "most major components of the literature lack an explicit theory of how the president, congress, bureaucracy, and courts interact" (Hammond and Knott 1996:120). There have been important theoretical insights on this issue: Fiorina (1981) posits that collective action problems among politicians are responsible for the poor monitoring of agencies, and Ferejohn and Shipman (1990) argue that the amount of policy consensus among politicians is a key determinant of successful implementation. However, it has proven difficult to develop these arguments further using formal modeling techniques. Because agency theory has no developed mechanism for handling multiple principals, few include this in their models (e.g., Wood 1988; Calvert, McCubbins, and Weingast 1989). The situation is similar to the way risk preferences are handled in economic agency theory: all note its theoretical importance, but most ignore it in practice to keep their mathematical models tractable. Perhaps agency theory works better for simpler state structures with unitary principals (like monarchies) than it does for complex contemporary states.

One of the most interesting omissions in the political science literature is any discussion of the recruitment or selection of agents. This issue is important in all other applications of agency theory—why not here? The main reason is probably the empirical narrowness of agency theory in political science. As in economics, agency theory in political science has been applied only to contemporary organizations in the United States—thus the questions asked in this literature reflect the issues that are important in this particular empirical setting. Most of the agents in bureaucratic agencies here are lifetime civil service employees, not chosen by current politicians, so agent selection has not been seen as a critical determinant of implementation outcomes. However, the issue of agent selection is not totally unimportant, even in the contemporary United States. The top-level agents in bureaucratic agencies are chosen by politicians, and their "type" might be important. If the historical scope of agency theory in political science were broader, it could be used to compare the patronage system for selection of state employees to the civil service system that followed it. Which of the two would "drift" less from politician's interests? Although we normally think of merit-based hiring as most efficient, one plausible alternative argument is that the patronage system would be better controlled since it selected for agents who shared the same values as the politicians who appointed them.[25] The outcome depends on the relative importance of technical competence and ideological loyalty.

Three things that are virtually ignored in economic agency models play an important role in political science arguments: third parties, administrative procedures, and multiple principals. State policy implementation always entails a triadic relationship among rulers of states, their officials, and citizens/subjects affected by policy. Agency theory formally treats only the first two, but there is ample empirical evidence that the third party is often equally important.[26] Political science has most fully incorporated third parties into agency theory. Administrative procedures were noted in one economics paper (Johnson and Libecap 1989), but their role too has only been fully explored in political science. The same is

---

[24] One move in this direction is found in a recent paper by Torenvlied (1996), the first to take the strength or salience of agents' preferences into account.

[25] Another possibility is that patronage ties could make agents dependent on particular politicians, thus limiting their corruption.

[26] Agency theory points to the potential importance of third parties, but contains no developed model to analyze them.

true of multiple principals—they are important in many contexts, but the political science literature is the first to fully assess their implications. In each of these cases, as agency theory moved to a new disciplinary context, new issues became significant. Of course, the political science literature has its blind spots, as well. The role of risk in agency relations is not discussed (perhaps because neither actor is ever a full residual claimant in the politician-bureaucrat relation). Although bureaucracies consist of multiple agents, agencies are generally treated as corporate actors, leaving their internal dynamics a "black box" and ignoring the possibly important consequences of collusion among bureaucratic agents. Finally, resource distributions do not play any role in most of these arguments. Wood (1988) is a notable exception, claiming that the greater the resources agents have, the more difficult they will be to control. Resources were also virtually ignored in the economic literature—sociologists using agency theory, beginning with Weber, stress this factor.

## THE CLASSICAL SOCIOLOGICAL ROOTS OF AGENCY THEORY: WEBER ON RULER-STAFF RELATIONS

For Weber, the relationship between rulers and their administrative staffs is essential to understanding political history. He stresses that rulers face a problem in controlling these agents, because the interests of the agents often differ from theirs.[27] He argues that "historical reality involves a continuous, though for the most part latent, conflict between chiefs and their administrative staffs for appropriation and expropriation in relation to one another" ([1922]1968:264). Furthermore, Weber realizes that the reason agents are often able to get away with acting contrary to the interests of principals is that they have better information concerning the quality of their performance than do principals (ibid.:225, 991–99). This is a classic statement of an agency problem, delegation of authority leading to problems of control due to conflicting interests of principals and agents, and informational asymmetries favoring agents. In general terms, Weber's arguments about how agency problems vary across administrative forms point to many of the same causal factors stressed by contemporary agency theorists—the type of agent chosen, the effectiveness of monitoring, and the nature of positive and negative sanctions.

A central component of Weber's analysis is microfoundations, since as a methodological individualist he believed that any complete explanation must include an account of the motives of actors (ibid.:4, 13).[28] Although Weber divided social action into four types, he argued that it would be *analytically useful* to begin all work by assuming instrumental microfoundations: "For the purposes of a typological scientific analysis it is convenient to treat all irrational, affectually determined elements of behavior as factors of deviation from a conceptually pure type of rational action" (ibid.:6).[29] Several discussions of Weber have noted that although his theoretical writings stress the existence of four different types of social action, many of his substantive analyses tend to be based primarily on instrumental microfoundations (Alexander 1985; Kahlberg 1994).[30]

Weber begins his discussion of tax administration by noting the multiple motives of officials. He suggests that officials obey rulers on the basis of some combination of their

---

[27] Weber also realized that the principals could engage in "appropriation and expropriation" of the resources of agents, a problem that Perrow (1990) correctly notes is ignored by contemporary agency theorists in economics. This is one of the many ways in which Weber's analysis is more substantively sophisticated, even if much less precise and formal, than contemporary economic agency theory.

[28] It is important to stress the difference between methodological individualism and reductionism, since Weber clearly rejected the latter. Individual action must always be part of a complete explanation, but will never be all of it.

[29] The fact that Weber was willing to begin his analysis with instrumental microfoundations does not mean that he thought instrumental rationality was ontologically primary—the choice was based on analytical utility.

[30] This is obviously not true of all of his substantive analyses—his sociology of religion, for example, is not based on instrumental microfoundations. Even his political sociology is not based solely on instrumental rationality, as noted below.

instrumental interests, custom, affectual ties, ideal interests, and legitimacy (ibid.:212–13). He then quickly narrows this long list. Weber recognized that the primary problem faced by rulers constructing systems of administration is that their officials do not totally share their interests (ibid.:264). He thought this problem could be mitigated in two main ways: (1) by constructing organizational forms that make it in the instrumental interests of officials to comply with the orders of rulers (by using appropriate forms of monitoring and sanctions); and (2) by fostering a belief in officials that the authority of the ruler is legitimate, and thus that it is their moral duty to obey. The microfoundations Weber used to model tax administration are thus based on a combination of instrumental rationality and value rationality, and the goals of actors are generally to maximize some combination of wealth[31] and power. I will focus on the models using instrumental rationality, and return to the more cultural issue of legitimacy at the end of this section.[32]

Weber's ideal types of forms of state organization (patrimonialism, bureaucracy) can best be understood as his attempt to model agency problems in different structural contexts. They are essentially typical clusters of recruitment, monitoring, and sanctioning strategies. Rulers (as principals) must delegate authority to some staff of state officials (as agents) in order to implement any of their policies. Weber argues that "[s]ociology seeks to formulate type concepts and generalized uniformities of empirical process. This distinguishes it from history, which is oriented to the causal analysis and explanation of individual actions, structures, and personalities possessing cultural significance" ([1922] 1968:19). He also has a clear preference for abstract over historically specific models: "The more sharply and precisely the ideal type has been constructed, thus the more abstract and unrealistic in this sense it is, the better it is able to perform its functions in formulating terminology, classifications, and hypotheses" (ibid.:21).[33] One of the important features of models is that they suggest which specific features of cases should be causally important. Even more fundamentally, the use of abstract ideal types allows for the separation of the general and particular features of each case, facilitating comparative analysis. Because of the variety of historical cases used to construct his ideal types, they cover a wide range of selection, monitoring, and sanctioning systems.

Weber was well aware that hiring agents of certain types would affect the extent of control problems in ruler-staff agency relations. For example, he argued that agency problems could be mitigated by hiring agents with high ability, using formal examination systems to identify merit (ibid.:264, 994, 1043). This form of recruitment is one cause of the efficiency of bureaucracy. Weber also situated this within a structural context, arguing that bureaucracy tends to be associated with a leveling of status differences (ibid.:225–26), in part because strong status distinctions inhibit merit-based hiring, and that it would not develop until educational institutions provided a sufficient number of trained potential officials (ibid.:973).[34] Recruitment based on merit is historically quite rare, so Weber focused primarily on selection based on prebureaucratic criteria.[35]

Weber also discussed the ways in which the process of agent selection affected the relative power of principals and agents (ibid.:232–33, 257, 1007). He was especially interested in the role of dependence in agency relations, since the more dependent the agents,

---

[31] Wealth is both easy to measure and is the most fungible of assets; it is necessary for rulers to satisfy a wide variety of idiosyncratic preferences, from the vast artistic treasures accumulated by Suleiman the Magnificent to Louis XIV's quest for glory through military victories.

[32] For an interesting attempt to model legitimacy from a rational choice perspective, see Levi's (1988, 1997) work on quasi-voluntary compliance and conditional cooperation.

[33] "In Weber's practiced methodology, 'sociology' is the generalized aspect of the study of history and contrasts with the causal analysis of individual phenomena" (Roth 1976:315).

[34] The causation goes both ways in Weber's argument—bureaucracy tends to further the process of status levelling.

[35] See his long list of patrimonial and extrapatrimonial forms of recruitment ([1922]1968:228–29).

the more likely they are to comply with principals' orders (ibid.:993, 1015–18, 1043). Hiring on the basis of high status (e.g., nobility) would decrease the dependence of agents (ibid.:1028), whereas hiring foreign or slave agents will increase dependence.[36] His focus on power and status considerations has unfortunately not influenced agency theorists outside sociology.

Weber was also interested in overall monitoring capacity and in exploring different organizational forms of monitoring. He focused most on the structural determinants of monitoring problems and the organizational forms that arose as attempts to solve them. His understanding of the importance of monitoring capacity is perhaps best indicated by his argument (ibid.:224) that the development of technologies of communications and transportation (essential foundations for adequate monitoring) were necessary conditions for the emergence of bureaucratic administration. Since monitoring problems increase with distance, the farther officials got from the ruler, the more they "evaded the ruler's influence" (ibid.:1051).

Weber discusses several ways to mitigate monitoring problems in prebureaucratic settings. For example, he argues that using collegiate organization (having multiple officials jointly fill one position in the hierarchy) and rotating officials can both facilitate monitoring. His discussion of the use of collegiate bodies also relies on the importance of information and the role of monitoring:

> This kind of collegiate body thus is the typical form in which the ruler, who increasingly turns into a dilettante, at the same time exploits expert knowledge and—what frequently remains unnoticed—seeks to fend off the threatening dominance of experts. He keeps one expert in check by others, and by such cumbersome procedures seeks personally to gain a comprehensive picture as well as the certainty that nobody prompts him into arbitrary decisions. (Weber [1922]1968:995; see also 222, 997).

Weber also realized that rotation could be used to mitigate monitoring problems, especially those due to collusion. By placing different agents in the same environment, rotation facilitates monitoring by helping rulers distinguish between outcomes due to the nature of the environment and those due to the actions of particular agents.[37]

Rulers also developed ways to decrease their monitoring problems by decreasing the number of units they taxed. One example of this is making communities collectively responsible for taxation: "The ruler would have required a very extensive coercive apparatus in order to get hold, in each instance, of the persons who were under liability, and this difficulty exactly was the reason for the system of compulsory associations upon which this task devolved" (ibid.:1024).[38]

Weber's discussion of the threat that bureaucratic administrative agencies pose to democracy in modern states also focuses on monitoring problems resulting from asymmetric information. Weber argued that elected politicians would find it increasingly difficult to monitor and control bureaucratic agents due the latter's possession of expert knowledge (ibid.:990–93). Further, Weber realized that officials would often find ways to prevent rulers from getting information about their actions, that their "striving for power" will often cause them to develop "official secrets" (ibid.:234–35; see also 992–93).

---

[36] He suggested that using local "notables" as agents would make administration slow and inefficient, because such agents had low dependence on rulers and strong network ties to other local elites ([1922]1968:974,1058–59, 1061–64).

[37] Rotation also decreases corruption by another mechanism—it limits the strength of ties between officials and taxpayers, thus mitigating problems of collusive corruption (Weber [1922]1968:1043).

[38] Weber refers to this as a system of "liturgical obligations" ([1922]1968:1022–25).

The use by rulers of various forms of positive and negative sanctions to control their agents was also central to Weber's account of state administration.[39] He explored the causes and the consequences of various forms of paying state agents, including payment in land and other benefices (ibid.:966, 1011, 1032),[40] status (ibid.:1028), and making agents full residual claimants in tax farming systems (ibid.:965). Furthermore, he compared the relative effectiveness of positive and negative sanctions, and of certainty versus severity of sanctions (ibid.:968). Weber always analyzes the use of sanctions within a structural context. Different sanctions require different amounts and types of resources, so resource distributions are key determinants of sanction forms. For example, "[s]ince the good will of officials depended on the possibility that their merits would be rewarded, the possession of a treasure was everywhere the indispensable basis of patrimonial domination" (ibid.:1038; see also 973). Weber goes on to argue that adequate resources are a necessary condition for the development of bureaucracy, since agents must regularly be paid fixed monetary salaries (ibid.:1059).

The use of the term "ruler" throughout this section masks another of Weber's important contributions, his discussion of the effects of variations in the form of rule (i.e., the nature of the principal) on the efficacy of administration. One of the best examples of this is his analysis of Roman tax farming ([1909]1976). Weber argues that Roman emperors more effectively controlled tax farmers (and administrators more generally) than the magistrates and senators who ruled the Roman Republic. Whereas "city-states, especially Rome, subjected their possessions to brutal exploitation by private capital through usurious tax farming" (ibid.:63), the "very first achievement of the Roman emperors had been the regulation of the tax system, when the arbitrary power of the tax farmers was curbed" (ibid.:364). The reason for the difference was the higher discount rates of the republican principals due to their short terms of rule. Emperors, with longer time horizons since they expect to rule for life, "aim at a prudent and durable rate of exploitation based on the actual resources and capacity of [their] subjects" (ibid.:63). Differences in the nature of principals, due to variations in the form of the state, thus have important impacts on administrative outcomes.

Weber was especially interested in the determinants of the degree of centralization of administration. The main issue in his discussion of patrimonialism is the problem of devolution—how the ruler can keep from totally losing control over such a decentralized form of administration. Weber's assumption throughout is that rulers would set up centralized administrative systems if they had the resources and technological capacity to do so, and that the lack of these is the main cause of decentralized patrimonialism ([1922]1968:1059).[41] Decentralization raises the risk of devolution: "fixed income in kind from the magazine of the lord or from his current intake—which has been the rule in Egypt and China for millennia and played and important part in the later Roman monarchy as well as elsewhere [has the disadvantage of moving] toward the appropriation of the sources of taxation by the official and their exploitation as private property" (ibid.:964).

Weber's discussion of feudalism is similar, since it is a subtype of patrimonialism. The main focus is on the necessity of decentralization and the power-driven problem of devolution (ibid.:257). For example, Weber argues that prebendal feudalism arises when rulers

---

[39] For example, see his long list of the possible forms of positive sanctions in patrimonial administration ([1922]1968:235–36).

[40] Weber notes that payment in benefices has a tendency to result in devolution, since rulers often lose control of agents they have made so independent ([1922]1968:1038–41).

[41] Weber argues that decentralization (in this case, decentralized patrimonialism) results from poor monitoring—the "weak development of the technical means of communication and therefore of political control" ([1922]1968:1091).

are unable to pay salaries to fund armies, and because it allows them to shift the risks associated with tax collection (ibid.:260–61).[42] The problem with feudalism is that "the latent struggle for authority becomes chronic between the lord and his vassal, and the ideal extent of feudal authority has never been effectively carried out in practice or remained effective on a permanent basis" (ibid.:257).

Weber's arguments about the causes of variations in forms of administration foreshadow those of contemporary agency theory in that they tend to stress efficacy or efficiency.[43] This is especially true of his analysis of the development of bureaucratic administration: "[t]he decisive reason for the advance of bureaucratic organization has always been its purely technical superiority over any other form of organization" (ibid.:973). Although Weber stresses efficiency, he also notes the importance of other power-based causes, both external and internal.[44] For example, he argues that war often affects the form of agency relations between rulers and state officials, basically by hastening the development of bureaucracy (ibid.:291, 966). However, war can also decrease administrative efficiency. For example, Weber sees the sale of offices as primarily caused by war: "the direct purchase of offices . . . occurs when the lord finds himself in a position in which he requires not only a current income but money capital—for instance, for warfare or debt payments" (ibid.:966). Moreover, Weber argues that internal power relations often prevented the use of the most efficient administrative forms. For example, "Whenever the prince could strengthen his position, his connections with all his subjects would become more direct in one way or another. However, as a rule the prince found himself compelled to compromise with the local patrimonial authorities or other honoratoires; he was restrained by the possibility of an often dangerous resistance, by the lack of a military and bureaucratic apparatus capable of taking over the administration and, above all, by the power position of the local honoratoires" (ibid.:1058).

Weber's model of the agency relation between rulers and their administrative staffs was not just an early formulation of the economic theory of agency—it was fundamentally sociological in the sense that it included both a rich depiction of the structural context within which the agency relation was embedded, and a complex view of the microfoundations of action. For Weber, ruler-staff agency relations are always embedded in a structural and historical context. Analyzing these relations thus requires not only the use of abstract ideal types (patrimonialism, bureaucracy, etc.), but concrete knowledge of historical conditions as well.

The most significant difference between Weber's approach to agency relations and that of most economists concerns the role of noninstrumental motivations, and at the macro

---

[42] This focus on who bears the risk in agency relationships will be central to economic agency theory; unfortuntely, it is not developed further by Weber.

[43] Weber's notion of the efficiency of bureaucracy is broader than that used by most contemporary economists, since it includes efficacy in broader political terms (legitimacy, security of rule) as well as the maximization of net tax revenue (although the latter was a major concern for Weber as well).

[44] Weber's discussion of tax farming illustrates the mix of efficiency and power as determinants of administrative forms. For example, "In the case of leasing [tax farming], the aim has been partly a practical financial one to meet the stringencies caused especially by the costs of war. It has partly also been a matter of the technique of financing, to insure a stable money income available for budgetary uses" ([1922]1968:234–35). In a later discussion (ibid.:965–66), Weber again stresses the importance of predictable revenue flows as a main determinant of tax farming, and goes on to analyze the causes of variations in the form of tax farming: "The lord seeks to safeguard himself against this loss of control by regulations. The mode of tax farming or the transfer of taxes can thus vary widely; depending upon the distribution of power between the lord and the farmer, the latter's interest in the full exploitation of the paying capacity of the subjects or the lord's interest in the conservation of this capacity may predominate. The nature of the tax farming system in the Ptolemaic empire, for instance, was clearly determined by the balance of the joint or the opposing influence of these motives: the elimination of oscillations in the yields, the possibility of budgeting, the safeguarding of the subjects' capacity to pay by protecting them against uneconomical exploitation, and state control of the tax farmer's yields for the sake of expropriating the maximum possible" (ibid.:965).

level, culture and legitimacy. [45] In addition to material factors, officials are more likely to comply with a ruler's directives if they believe the ruler is legitimate—by which he means a belief that the ruler has the *right* to give orders and officials have a *duty* to obey them. The extent to which officials will comply with rulers' commands, and thus the effectiveness of the administrative system, is thus for Weber a joint function of material incentives and the extent and type of legitimacy. His three types of legitimate domination basically specify different contexts within which agency relations between rulers and officials are embedded, combining features of organizational structure and culture.

Weber's inclusion of complex microfoundations and culture makes it more difficult to construct and test clear propositions, but may also allow agency theory to be applied to a wider range of cases, and to explain a greater proportion of the variance in outcomes. Contemporary sociologists using agency theory have chosen to emphasize different aspects of Weber's arguments. Some rely mainly or exclusively on organizational structure and incentives, the side of Weber most consistent with rational choice (Kiser, 1994; Kiser and Schneider 1994; Adams 1996), whereas others incorporate legitimacy and culture into their analyses of agency relations, thus going beyond the bounds of traditional rational choice theory (Hamilton and Biggart 1980, 1984; Gorski 1993).[46]

## AGENCY THEORY IN SOCIOLOGICAL ANALYSES OF STATE POLICY IMPLEMENTATION

Just as the firm was generally viewed as a "black box" in economic analysis, the state was usually modeled as a unitary corporate actor in most sociological theories. This has begun to change only recently. Compared to both economics and political science, the use of agency theory in sociology is in its infancy. It does provide an interesting comparison, however, since its intellectual genealogy is quite different. In contrast to political scientists, sociologists using agency theory to study state policy implementation have drawn on both Weberian and economic versions of agency theory. This difference in part reflects disciplinary boundaries, but is also due to the fact that agency theory in political sociology has been applied primarily to early modern states—a context in which Weber's seminal arguments cannot be ignored (whereas both economists and political scientists have focused on the contemporary United States). This illustrates a general point: the empirical context to which a theory is applied often has important effects on its content.

Neo-Weberian agency analyses of state administration have both elaborated on and challenged many of Weber's substantive conclusions. Four recent examples are Adams's (1996) discussion of agency problems involved in colonial control, Kiser's (Kiser and Tong 1992; Kiser and Schneider 1994; Kiser 1994) work on the organizational structure of early modern tax administration, Gorski's (1993) analysis of the "disciplinary revolution" in Holland and Prussia, and Hamilton and Biggart's (1980, 1984, 1985) arguments about control in contemporary state government bureaucracies. Since both Adams and Kiser identify their work as Weberian and explicitly use rational choice and agency theory, I

---

[45] For example, Weber's substantive comments about compliance within feudal organizations rest primarily on instrumental microfoundations, but other motives are mentioned as well. Weber suggests that both "honor" and a "solidary, fraternal relationship" contribute to the vassal's compliance with the demands made by the lord ([1922]1968:1077–78, 255). However, he argues that such "voluntary obedience" and "purely personal loyalty" are never sufficient, and thus that in the absence of material incentives the "lord's authority is precarious" (ibid.:257).

[46] Even the sociologists who incorporate a broader range of microfoundations still focus to a large extent on instrumental motives, organizational structure, and material incentives. Their work could thus be understood as adding to the baseline rational choice model of agency, not totally rejecting it (although neither Gorski nor Hamilton would describe their work this way). This raises the issue of the definition and scope of rational choice theory, unfortunately much too large a topic to address here.

begin with a discussion of their contributions and conclude with an analysis of the work of Hamilton and Biggart and Gorski. These scholars draw on Weber in rather different ways. Both Adams and Kiser focus predominantly on the materialist and rational choice aspects of Weber's arguments, whereas Hamilton and Biggart and Gorski incorporate culture and heterogeneous microfoundations to a larger extent.

Adams's (1996) deft analysis of the agency relationship between patrimonial states and colonial trading companies in the Netherlands and England illustrates the affinity between sociological versions of agency theory and a Weberian approach. She begins with standard rational choice microfoundations,[47] then uses agency theory to model the relationship between metropolitan principals and colonial "company men." Her use of agency theory draws on economics, but is fundamentally Weberian. The agency model in Adams's work is not sparse and ahistorical but is deeply embedded in particular structural conditions. For example, both the Dutch and English states were patrimonial, so she situates her analysis within the central dynamic of patrimonialism identified by Weber, the tendency of principals to totally lose control of agents, resulting in devolution.[48] Both patrimonial states also suffered from multiple principals (what Adams calls the "hydra factor")—"multiple heads or principals lacked institutional mechanisms to resolve the resulting uncertainties and infighting" (1996:16). As political scientists (Banfield 1975; Fiorina 1981; Hammond and Knott 1996) have noted, the existence of multiple principals clearly magnifies problems of control.

Given the problems principals in the Netherlands and England faced gathering information about the activities of their agents in Asia due to the distance involved, how were they able to control them? Her central focus is on problems of monitoring and control, and how they were influenced by variations in organizational structure. Adams uses network theory to model aspects of organizational structure, and thus to explain variations in levels and types of corruption in the administration of trading companies. In the Dutch hierarchy, the Batavian outpost (contemporary Jakarta) maintained a middleman or brokerage position between the metropolitan principals and company agents, which allowed them to illegally extract some of the surplus by collusive corruption with other agents. In spite of this, the level of corruption was limited. To account for this, Adams makes an argument that was stressed often by Weber ([1922]1968:1007, 1015–18), that the level of corruption is inversely related to the level of agent dependence on principals. Agents of the Dutch had no alternative opportunities, and this dependence limited their corruption. The situation changed when the English company moved in, since it provided not only direct competition, but also more opportunities for collusive corruption for agents of the Dutch (Adams 1996:23). The ultimate result is the one predicted by Weber: by the end of the eighteenth century patrimonial principals had almost totally lost control of their colonial agents.

Kiser also uses agency theory to address classic Weberian questions. For example, Kiser and Tong (1992) explore the levels and types of corruption in Ming and Qing tax administration in part to discover whether instrumental motivations were dominant in this premodern, nonwestern setting. They demonstrate that tax agents in premodern Asia reacted to variations in monitoring and sanctioning (the certainty and severity of punishments for corruption) just as their western counterparts did. Moreover, in contrast to cultural social-

---

[47] Adams assumes that "both principals and agents tend to act in intendedly rational fashion, and opportunistically, to advance their own individual gains" (1996:14). In the conclusion of the paper, Adams (p. 26) broadens the microfoundations of her argument by briefly exploring the possibility that motivations such as family honor and position may have been important, as well. This too illustrates the statregy advocated by Weber and here—that scholars should begin with simple instrumental microfoundations and then expand them if necessary.

[48] A second consequence of the patrimonial context is a low level of division of labor. As Adams puts it, "the substantive content of key roles and ties is multivocal" (1996:15). Key agents had multiple goals, both economic and coercive.

ization arguments, they demonstrate that the Confucian education and examination system did not reduce corruption by giving officials prostate values, but in fact increased it by creating strong long-term network ties that facilitated collusion among officials.

Kiser and Schneider (1994) address Weber's claim that the efficiency of Prussian tax administration was due to its early bureaucratization. Using historical data not available in Weber's time, they show that the Prussian state was much less bureaucratic than Weber thought. Moreover, they demonstrate that particular variations from the bureaucratic ideal type that increased the dependence of agents or strengthened their incentives were the primary causes of efficiency in this case. For example, Prussian rulers used a unique system of caring for injured military veterans. Instead of welfare payments, they gave them positions as collectors of indirect taxes (what we would now call a "workfare" program). Since these officials had poor alternative employment opportunities, they were very dependent on rulers, and thus less corrupt. By creating a high level of dependence, this way of selecting officials was more effective than bureaucratic selection on the basis of merit.

Kiser's (1994) comparative study of early modern tax farming also revises Weber's conclusions. Weber ([1922]1968:965) argued that tax farming was less efficient than state administration, and that it was only used to make revenue flows more predictable. By looking at tax farming and bureaucratic fixed salaries as alternative forms of agency relations, Kiser demonstrates that tax farming was more efficient than state administration for the collection of indirect taxes, since it not only provided stronger incentives for agents but also mitigated specific monitoring problems.

Gorski (1993) has also drawn on Weber to develop an argument about tax administration in the early modern era,[49] by applying to the state Weber's main thesis in *The Protestant Ethic and the Spirit of Capitalism* ([1920]1996).[50] Part of his argument implicitly uses a version of agency theory. Throughout the paper, Gorski (1993:271–72) focuses on the ways in which various "material interests," "externalized controls," and "institutional restraints" limited corruption and thus reduced administrative costs in Holland and Prussia. He concentrates on changes in the form and intensity of monitoring techniques, such as the use of ad hoc commissioners and "colleges" (ibid.:287, 297).

Gorski thinks instrumental motivations and agency theory cannot fully explain the efficiency of state administration in Holland and Prussia. He also focuses on the role of internalized religious values and religious monitoring mechanisms in motivating compliance by tax officials. He argues that one of the main causes of the efficiency of administration in Prussia and Holland was that rulers selected agents on the basis of religious affiliation, and that these agents had religious values that inhibited corruption. This would normally be coded as a cultural argument, the antithesis of a rational choice account using agency theory. However, what Gorski is doing theoretically can easily be put in the language of agency theory: he concentrates on the consequences of hiring agents of different types, differentiated by religious affiliation (which he takes as indicating internalized religious values). Gorski has added to these economic models by making agent type endogenous. Instead of simply stipulating that agent "type" is given by "nature," he argues that it is derived from religious affiliation.

Economic agency theorists might suggest a different causal mechanism to account for the correlation between religious affiliation and agent compliance.[51] Drawing on Spence

---

[49] Gorski's paper is much broader than this, dealing with issues of state formation as well, but we focus only on administrative issues.

[50] See Coleman (1986) on the affinity between Weber's "Protestant ethic" thesis and rational choice theory.

[51] Since I am interested in the form of the theoretical argument and not whether or not it is empirically correct, I will not address the question of whether the correlation in fact exists. For a debate about this, see Gorski (1995) and Kiser and Schneider (1995).

(1973), one could argue that joining a religious organization opposed to corruption is a signaling device that potential state officials could use to indicate honesty. If churches did punish corrupt behavior, this would be an effective signal since it would be cheaper for honest than for dishonest officials to join the church. The causal relation is the same (religious affiliation would be correlated with honesty), but the causal mechanism in the signaling argument does not involve the internalization of religious values.

Gorski (1993) also suggests an additional causal mechanism linking religious affiliation with low corruption, third-party monitoring and sanctioning by churches. If corruption is contrary to religious as well as state rules, and churches enforce those rules, they are indirectly acting as control mechanisms for the state. This part of Gorski's argument not only does not contradict agency theory, but provides additional evidence for the large literature in political science on third-party monitoring (Weingast 1984; Kiewiet and McCubbins 1991).

Hamilton and Biggart's (1980, 1984, 1985) work on policy implementation concentrates on Weber's central question: how can rulers control bureaucrats, given that the ruler is a dilettante and the bureaucrats are experts? They argue (1980) that in both early modern and contemporary states, leaders have tried to solve this problem by using a patrimonial intermediary, in the form of "personal staffs," to help them control bureaucrats. The personal staff is essentially a personalistic device for monitoring bureaucratic agents. In more general terms, they come to the interesting conclusion that the control of bureaucracies necessarily requires nonbureaucratic elements.

This line of argument is developed in much more detail in a fascinating comparison of the mechanisms of control used by two governors of California, Ronald Reagan and Jerry Brown (Hamilton and Biggart 1984). Although they make some interesting arguments about the differences between the two (mainly focusing on differences in the amount and type of delegation—as one would guess, Reagan delegated more), most of the analysis concentrates on the similarities in the agency problems they faced. Their argument mainly looks at instrumental interests (see Hamilton and Biggart [1985:15–16] for an explicit discussion of their primacy), but also incorporates values. Hamilton and Biggart (1984) stress that different types of agents are controlled in different ways. The personal staff is selected on the basis of personal ties and loyalty; they are usually dependent on the governor, and if all else fails they can be sanctioned severely and arbitrarily—in Weber's terms, the relationship is patrimonial (1985:29–33). Compliance of cabinet heads is maintained less by personal ties than by philosophical and ideological similarity (ibid.:55–66). This argument draws on Weber's stress on the importance of legitimacy, and is in many ways similar to Gorski's account of the "disciplinary revolution." Since cabinet heads have a great deal of autonomy from the governor, monitoring is difficult, so a similarity of fundamental preferences is important to ensure that they usually act in the interests of the governor. Of course, when they act contrary to the governor's interests, sanctions (such as cutting their budgets) are used as well (ibid.:95–97). Finally, they argue that professional experts are controlled by a combination of monitoring and sanctions typically used for civil servants and by drawing on their professional loyalty (in a manner not dissimilar to the way the church controls officials in Gorski's account).

Sociological versions of agency theory have used the heritage of Weber and some core ideas from their discipline to expand agency theories of policy implementation in useful ways. All of them have followed Weber in maintaining a broad empirical scope, and in concentrating on the structural contexts within which agency relations take place. Features of social structure that for economists are exogenously given by "nature" are often the focus of sociological analyses. Sociologists have begun to address questions of multiple agents and multiple principals, and have employed models from network theory to develop more formal and rigorous analyses of their interactions (Adams 1996). They have also

concentrated on resource distributions in agency relations, especially how they affect the dependence of agents (Kiser and Schneider 1994; Adams 1996). Finally, they have recently tried to address the cultural determinants of variations in agent types and agent compliance (Hamilton and Biggart 1980, 1984; Gorski 1993), a necessary foundation for a more complete theory of agent selection.

There are interesting differences in the use of agency theory among contemporary sociologists. Although all focus on material constraints and instrumental incentives (and are in fact much more similar than some rhetorical confrontations would suggest), the addition of cultural factors and heterogeneous microfoundations by Hamilton and Biggart and by Gorski makes their arguments substantively richer and more complete, but also more difficult to test.

CONCLUSION

This review demonstrates that the process of the intellectual diffusion of agency theory from economics to other disciplines is not well characterized by the polemical term "economic imperialism." Agency theory has not only transformed work on organizations in several disciplines, it has been transformed by these different disciplinary contexts. This is clearly not a case of "economic imperialism," but one of selective borrowing shaped by particular intellectual genealogies and disciplinary norms. Political scientists have drawn on economic models instead of Weber, and have in several respects remained fairly close to economic models of agency (retaining their parsimony at both micro and macro levels). In contrast, agency theory in sociology has been based on much broader conceptions of both the micro and macro levels—developing along lines very different from economic agency theory. This is a product of its intellectual genealogy (coming mainly from Weber), and of the core ideas and modes of explanation currently dominant in the discipline. Intellectual diffusion in this case (and probably most others) is not best described as imperialism, but as a complex form of assimilation, in which ideas from other disciplines both shape and are shaped by their new disciplinary contexts.

Perhaps even more striking is the fact that agency theory has also transformed its home discipline. Agency theory is part of the "new institutionalism" that is moving beyond the sparse mathematical models that have dominated neoclassical economics, and attempting to draw a more complete and nuanced picture of organizations and exchanges. As agency theory has evolved in economics, and especially as it has been applied to the study of political organizations, it has begun to include more complexity at both micro and macro levels of analysis. Some economic agency theorists are now incorporating aspects of social and organizational structure, multiple agents, and even heterogeneous microfoundations.[52]

What can we now conclude about the current value of agency theory and its future prospects? Contemporary opinions about agency theory are diverse. Jensen refers to it as a "revolution in the science of organizations" that provides all of the necessary "major analytical building blocks of a theory of organizations" (1983:321). Petersen sees much of value in economic agency theory, and wants to bring it "into the sociological mainstream" in industrial and organizational sociology (1993:227). On the other hand, Perrow argues that agency theory is "not only wrong but dangerous" because it reflects conservative ideology, and that although "we all invoke agency theory," we do so "in our worst moments" (1990:121,127).

It should be clear from this review of the varieties of agency theory that Perrow's polemical criticisms are misguided. Although economic agency theory is limited in several

---

[52] There is good reason to believe that the broadening of economic agency models will continue. The normal form in the discipline is for work in an area to begin with very sparse models and many strict assumptions. Later work in the area usually elaborates on the models and relaxes some of the more unrealistic assumptions. Agency theory is still quite new in economics; it is just now entering this second stage.

respects, his appraisal is far too negative even for the economic version (which is evolving in many ways to address his criticisms), and most of his criticisms do not apply at all to the much broader work using agency theory in other disciplines. Sociologists should not be taking a protectionist stance toward foreign imports—it is becoming increasingly clear that free trade in ideas benefits all parties, since most intellectual progress is taking place in the intersections of disciplines. Of course, it is also far too soon to conclude that agency theory has produced the "revolution" in the study of organizations hailed by Jensen (1983:321). Although I hope to have demonstrated that agency theory provides a promising model of state administration, it would take a more systematic comparative analysis to prove that it is superior to other existing theories.

In spite of all of the recent talk about the importance of interdisciplinary exchanges, the case of agency theory suggests that we have not moved very far in that direction. Although the general model has moved across disciplines, the particular insights and arguments made in one discipline are often ignored by others. We seem to be interdisciplinary in a fairly superficial way. The future prospects of agency theory in large part depend on breaking down the disciplinary barriers that up until now have kept these literatures fairly isolated from each other. There is much that scholars in each discipline could learn from the others. Economists could learn from Weber and contemporary sociologists and political scientists how to incorporate social structure, organizational forms, multiple principals, social networks, and third parties into their models. Much of this will follow if they simply expand the empirical scope of their work. The same is true for political science—they could learn from economists' work on risk, and from sociological analyses of social and organizational structure and networks. Sociologists could also benefit from incorporating risk, paying more attention to problems of multiple principals, and trying to increase the rigor and precision of their theoretical models.

The future of agency theory in the study of state policy implementation is open and uncertain, since it is currently developing in several different directions. Analyses using agency theory range from precise but narrow economic models to neo-Weberian historical analyses. The same can be said for rational choice theory more generally. Standard criticisms of rational choice that equate it with neoclassical economics can no longer be taken very seriously. Rational choice in sociology is clearly not the same as neoclassical economics, since it often incorporates both broader microfoundations and richer models of social structure (Kiser and Hechter 1998).[53] It is a different type of theory, just now beginning to take shape.

# REFERENCES

Adams, Julia. 1996. "Principals and Agents, Colonialists and Company Men: The Decay of Colonial Control in the Dutch East Indies." *American Sociological Review* 61(Feb):12–28.

Alexander, Jeffrey. 1985. *Theoretical Logic in Sociology, Vol 3: The Classical Attempt at Theoretical Synthesis—Max Weber*. Berkeley: University of California Press.

Arrow, Kenneth. 1951. *Social Choice and Individual Values*. New York: John Wiley and Sons.

———. 1964. *Essays in the Theory of Risk-Bearing*. Chicago: Aldine.

Bac, Mehmet. 1996. "Corruption, Supervision, and the Structure of Hierarchies." *Journal of Law, Economics, and Organization* 12(2):277–99.

Banfield, Edward. 1975. "Corruption as a Feature of Governmental Organization." *Journal of Law and Economics* 18:587–605.

[53] Varieties of sociological rational choice theory have been used to reanalyze many classical sociological issues, including the causes of group solidarity (Hechter 1987), the emergence of norms (Coleman 1990), variations in gender stratification (Brinton 1988), the rise of fascism (Brustein 1996), state tax policies (Levi 1988), the causes of war initiation (Kiser, Drass, and Brustein 1995), and the relationship between class struggle and racial conflict (Brown and Boswell 1995), to list just a few.

Barzel, Yoram. 1989. *Economic Analysis of Property Rights*. Cambridge: Cambridge University Press.

Bates, Laurie, and Rexford Santerre. 1993. "Property Tax Collector Performance and Pay." *National Tax Journal* 45(1):23–30.

Bawn, Kathleen. 1995. "Political Control Versus Expertise: Congressional Choices about Administrative Proce-dures." *American Political Science Review* 89(1):62–73.

Becker, Gary. 1957. *The Economics of Discrimination*. Chicago: University of Chicago Press.

———. 1976. *The Economic Approach to Human Behavior*. Chicago: University of Chicago Press.

Becker, Gary, and George Stigler. 1974. "Law Enforcement, Malfeasance, and Compensation of Enforcers." *Journal of Legal Studies* 3:1–18.

Bendor, Jonathan, and Terry Moe. 1985. "An Adaptive Model of Bureaucratic Politics." *American Political Science Review* 79:755–74.

Berhold, Marvin. 1971. "A Theory of Linear Profit-Sharing Incentives." *Quarterly Journal of Economics* 85(3):460–82.

Bernstein, M. 1955. *Regulating Business by Independent Commission*. Princeton: Princeton University Press.

Border, K.C., and J. Sobel. 1987. "Samurai Accountant: A Theory of Auditing and Plunder." *Review of Economic Studies* 54:524–40.

Brinton, Mary. 1988. "The Social-Institutional Bases of Gender Stratification: Japan as an Illustrative Case." *American Journal of Sociology* 94(2):300–34.

Brown, Cliff, and Terry Boswell. 1995. "Strikebreaking or Solidarity in the Great Steel Strike of 1919: A Split Labor Market, Game-Theoretic, and QCA Analysis." *American Journal of Sociology* 100(6):1479–1519.

Brustein, William. 1996. *The Logic of Evil*. New Haven: Yale University Press.

Buchanan, James, and Gordon Tullock. 1962. *The Calculus of Consent*. Ann Arbor: University of Michigan Press.

Calvert, Randall, Matthew McCubbins, and Barry Weingast. 1989. "A Theory of Political Control and Agency Discretion." *American Journal of Political Science* 33:588–61.

Chandler, P., and L. Wilde. 1992. "A General Characterization of Optimal Income Taxation and Enforcement." Social Science Working Paper 791, California Institute of Technology, Pasadena.

Cheung, Stephen. 1969. *A Theory of Share Tenancy*. Chicago: University of Chicago Press.

Coase, Ronald. 1937. "The Nature of the Firm." *Economica* 4:386–405.

Coleman, James. 1986. "Social Theory, Social Research, and a Theory of Action." *American Journal of Sociol-ogy* 91(6):1309–35.

———. 1990. *Foundations of Social Theory*. Harvard: Belknap Press.

Cremer, H., M. Marchand, and P. Pestieau. 1990. "Evading, Auditing, and Taxing: The Equity-Compliance Trade-off." *Journal of Public Economics* 43:67–92.

Demski, Joel, and David Sappington. 1987. "Hierarchical Regulatory Control." *RAND Journal of Economics* 18(3):369–83.

Downs, Anthony. 1957. *An Economic Theory of Democracy*. New York: Harper and Row.

———. 1967. *Inside Bureaucracy*. Boston: Little, Brown.

Ehrenberg, Ronald, Richard Chaykowski, and Randy Ehrenberg. 1988. "Determinants of the Compensation and Mobility of School Superintendents." *Industrial and Labor Relations Review* 41(April):386–401.

Eisenhardt, K. 1989. "Agency Theory: An Assessment and Review." *Academy of Management Review* 14(1):57–74.

Ferejohn, John, and Charles Shipman. 1990. "Congressional Influence on Bureaucracy." *Journal of Law, Eco-nomics, and Organization* 6:S1–S20.

Fiorina, Morris. 1986. "Congressional Control of the Bureaucracy: A Mismatch of Incentives and Capabilities" in Lawrence C. Dodd and Bruce I. Oppenheimer (eds) *Congress Reconsidered*. Washington: Congressional Quarterly Press.

Fiorina, Morris. 1990. "Comment: The Problems with PPT." *Journal of Law, Economics, and Organization* 6:255–61.

Freeman, J.L. 1955. The Political Process: Executive Bureau-Legislative Committee Relations. New York: Ran-dom House.

Friedman, Debra, and Michael Hechter. 1988. "The Contribution of Rational Choice Theory to Macrosociolog-ical Research." *Sociological Theory* 6:201–18.

Goldstein, Gerald, and Ronald Ehrenberg. 1976. "Executive Compensation in Municipalities." *Southern Eco-nomic Journal* 43(July):937–47.

Gorski, Philip. 1993. "The Protestant Ethic Revisited: Disciplinary Revolution and State Formation in Holland and Prussia." *American Journal of Sociology* 99(2):265–316.

Gorski, Philip S. 1995. "The Protestant Ethic and the Spirit of Bureaucracy" *American Sociological Review* 60(5):783–86.

Graetz, M.J., J. Reinganum, and L. Wilde. 1986. "The Tax Compliance Game: Toward an Interactive Theory of Law Enforcement." *Journal of Law, Economics, and Organization* 2:1–32.

Greif, Avner. 1994. "Cultural Beliefs and the Organization of Society: A Historical and Theoretical Reflection on Collectivist and Individualist Societies." *Journal of Political Economy* 102:912–50.

Hamilton, Gary, and Nicole Woolsey Biggart. 1980. "Making the Dilettante an Expert: Personal Staffs in Public Bureaucracies." *Journal of Applied Behavioral Science* 16:192–210.

———. 1984. *Governor Reagan, Governor Brown*. New York: Columbia University Press.

———. 1985. "Why People Obey: Theoretical Observations on Power and Obedience in Complex Organizations" *Sociological Perspectives* 28(1):3–28.

Hammond, Thomas, and Jack Knott. 1996. "Who Controls the Bureaucracy?: Presidential Power, Congressional Dominance, Legal Constraints, and Bureaucratic Autonomy in a Model of Multi-Institutional Policy-Making." *Journal of Law, Economics, and Organization* 12(1):119–66.

Hechter, Michael. 1987. *Principles of Group Solidarity*. Berkeley: University of California Press.

Hechter, Michael and Satoshi Kanazawa. 1997. "Sociological Rational Choice Theory" *Annual Review of Sociology* 23:191–214.

Heckathorn, Douglas. 1988. "Collective Sanctions and the Creation of Prisoner's Dilemma Norms." *American Journal of Sociology* 94:535–62.

Huntington, Samuel. 1952. "The Marasmus of the ICC: The Commission, the Railroads, and the Public Interest." *Yale Law Journal* 614:467–509.

Jensen, Michael. 1983. "Organizational Theory and Methodology." *Accounting Review* 58(2):321–39.

Jensen, Michael, and William Meckling. 1976. "Theory of the Firm: Managerial Behavior, Agency Costs, and Ownership Structure." *Journal of Financial Economics* 3:305–60.

Johnson, Ronald, and Gary Libecap. 1989. "Bureaucratic Rules, Supervisor Behavior, and the Effect on Salaries in the Federal Government." *Journal of Law, Economics, and Organization* 5(1):53–81.

Kalberg, Stephen. 1994. *Max Weber's Comparative-Historical Sociology*. Chicago: University of Chicago Press.

Kiewiet, D. Roderick, and Mathew McCubbins. 1991. *The Logic of Delegation*. Chicago: University of Chicago Press.

Kiser, Edgar. 1989. "A Principal-Agent Analysis of the Initiation of War in Absolutist States." Pp. 65–82 in *War in the World System*, edited by Robert Schaeffer. New York: Greenwood.

———. 1994. "Markets and Hierarchies in Early Modern Fiscal Systems: A Principal-Agent Analysis." *Politics and Society* 22(3):284–315.

Kiser, Edgar, and Kathryn Baker. 1994. "Could Privatization Increase the Efficiency of Tax Collection in Less Developed Countries?" *Policy Studies Journal* 22(3):489–500.

Kiser, Edgar, and Michael Hechter. 1998. "The Debate on Historical Sociology: Rational Choice Theory and its Critics" *American Journal of Sociology* 104(3):785–816.

Kiser, Edgar, and Xiaoxi Tong. 1992. "Determinants of the Amount and Type of Corruption in State Fiscal Bureaucracies: An Analysis of Late Imperial China." *Comparative Political Studies* 25:300–1.

Kiser, Edgar, and Joachim Schneider. 1994. "Bureaucracy and Efficiency: An Analysis of Taxation in Early Modern Prussia." *American Sociological Review* 59(April):187–204.

Kiser, Edgar, and Joachim Schneider. 1995. "Rational Choice Versus Cultural Explanations of the Efficiency of the Prussian Tax System" *American Sociological Review* 60(5):787–91.

Kiser, Edgar, Kriss A. Drass, and William Brustein. 1995. "Ruler Autonomy and War in Early Modern Western Europe." *International Studies Quarterly* 39:109–38.

Klitgaard, Robert. 1988. *Controlling Corruption*. Berkeley: University of California Press.

Kofman, Fred, and Jacques Lawarree. 1996. "A Prisoner's Dilemma Model of Collusion Deterrence." *Journal of Public Economics* 59:117–36.

Kotowitz, Y. 1987. "Moral Hazard." Pp. 207–13 in *Allocation, Information, and Markets*, edited by John Eatwell, Murray Milgate, and Peter Newman. New York: Norton.

Levi, Margaret. 1988. *Of Rule and Revenue*. Berkeley: University of California Press.

———. 1997. *The Contingencies of Consent*. Cambridge: Cambridge University Press.

MacDonald, Glenn. 1984. "New Directions in the Economic Theory of Agency." *Canadian Journal of Economics* 37(3):415–40.

McCubbins, Mathew. 1985. "The Legislative Design of Regulatory Structure." *American Journal of Political Science* 29:721–48.

McCubbins, Mathew, and Thomas Schwartz. 1984. "Congressional Oversight Overlooked: Policy Patrols vs, Fire Alarms." *American Journal of Political Science* 28:165–79.

McCubbins, Mathew, Roger Noll, and Barry Weingast. 1987. "Administrative Procedures as Instruments of Political Control." *Journal of Law, Economics, and Organization* 3(2):243–77.

Michels, Robert. [1915]1959. *Political Parties*. New York: Dover.

Moe, Terry. 1985. "Control and Feedback in Economic Regulation: The Case of the NLRB." *American Political Science Review* 79:1094–116.

———. 1990. "Political Institutions: The Neglected Side of the Story." *Journal of Law, Economics, and Organization* 6:213–53.

Mookherjee, D., and I P'ng. 1989. "Optimal Auditing, Insurance, and Redistribution." *Quarterly Journal of Economics* 103:399–415.

Niskanen, William. 1971. *Bureaucracy and Representative Government.* Chicago: Aldine-Atherton.

North, Douglass. 1981. *Structure and Change in Economic History*. New York: Norton.

North, Douglass, and Robert Thomas. 1973. *The Rise of The Western World*. Cambridge: Cambridge University Press.

Peltzman, Sam. 1976. "Toward a More General Theory of Regulation." *Journal of Law and Economics* 19:211–40.

Perrow, Charles. 1990. "Economic Theories of Organization" Pp. 121–152 in *Structures of Capital*, edited by Sharon Zukin and Paul Dimaggio. Cambridge: Cambridge University Press.

Petersen, Trond. 1993. "The Economics of Organization: The Principal-Agent Relationship." *Acta Sociologica* 36:277–93.

Reinganum, J.R., and L.L. Wilde. 1985. "Income Tax Compliance in a Principal-Agent Framework." *Journal of Public Economics* 26:1–18.

Rose-Ackerman. 1978. *Corruption: A Study in Political Economy*. New York: Academic Press.

Ross, Stephen. 1973. "The Economic Theory of Agency: The Principal's Problem." *American Economic Review* 63(2):134–39.

Roth, Gunther. 1976. "History and Sociology in the Work of Max Weber." *British Journal of Sociology* 27(3):306–17.

Sanchez, Isabel, and Joel Sobel. 1993. "Hierarchical Design and Enforcement of Income Tax Policies." *Journal of Public Economics* 50:345–69.

Schelling, Thomas. 1978. *Micromotives and Macrobehavior*. New York: Norton.

Scotchmer, S. 1986. "Equity and Tax Enforcement." Harvard Discussion Paper.

Shliefer, Andrei, and Robert Vishny. 1993. "Corruption." *Quarterly Journal of Economics* August:599–617.

Smith, Adam. [1776]1976. *The Wealth of Nations*. Chicago: University of Chicago Press.

Spence, A.M. 1973. *Market Signalling: Information Transfer in Hiring and Related Processes*. Cambridge: Harvard University Press.

Stigler, George. 1971. "The Theory of Economic Regulation." *Bell Journal of Economics* 2:3–21.

Stiglitz, Joseph. 1987. "Principal and Agent." Pp. 241–53 in *Allocation, Information, and Markets*, edited by John Eatwell, Murray Milgate, and Peter Newman. New York: Norton.

Tirole, Jean. 1986. "Hierarchies and Bureaucracies: On the Role of Collusion in Organizations." *Journal of Law, Economics, and Organization* 2(2):181–214.

———. 1992. "Collusion and the Theory of Organizations." Pp. 151–206 in *Advances on Economic Theory: Sixth World Congress*, Vol. 2, edited by Jean-Jacques Laffont. Cambridge: Cambridge University Press.

Torenvlied, Rene. 1996. "Political Control of Implementation Agencies: Effects of Political Consensus on Agency Compliance." *Rationality and Society* 8(1):25–56.

Truman, David. 1951. *Governmental Process: Political Interests and Public Opinion*. New York: Knopf.

Weber, Max. [1909]1976. *The Agrarian Sociology of Ancient Civilizations*. London: New Left Books.

———. [1920]1996. *The Protestant Ethic and the Spirit of Capitalism*. Los Angeles: Roxbury.

———. [1922]1968. *Economy and Society*. Berkeley: University of California Press.

Weingast, Barry. 1984. "The Congressional-Bureaucratic System: A Principal-Agent Perspective." *Public Choice* 44:147–92.

Weingast, Barry, and M. Moran. 1983. "Bureaucratic Discretion or Congressional Control?: Regulatory Policy Making by the Federal Trade Commission." *Journal of Political Economy* 91:765–800.

Williamson, Oliver. 1976. *Markets and Hierarchies*. New York: Free Press.

———. 1985. *Economic Insititutions of Capitalism*. New York: Free Press.

Wood, B. Dan. 1988. "Principals, Bureaucrats, and Responsiveness in Clean Air Enforcements." *American Political Science Review* 82(1):213–234.

# Agency Problems and Dividend Policies around the World

RAFAEL LA PORTA, FLORENCIO LOPEZ-DE-SILANES,
ANDREI SHLEIFER, and ROBERT W. VISHNY*

## ABSTRACT

This paper outlines and tests two agency models of dividends. According to the "outcome model," dividends are paid because minority shareholders pressure corporate insiders to disgorge cash. According to the "substitute model," insiders interested in issuing equity in the future pay dividends to establish a reputation for decent treatment of minority shareholders. The first model predicts that stronger minority shareholder rights should be associated with higher dividend payouts; the second model predicts the opposite. Tests on a cross section of 4,000 companies from 33 countries with different levels of minority shareholder rights support the outcome agency model of dividends.

THE SO-CALLED DIVIDEND PUZZLE (Black (1976)) has preoccupied the attention of financial economists at least since Modigliani and Miller's seminal work (see Modigliani and Miller (1958) and Miller and Modigliani (1961)). This work established that, in a frictionless world, when the investment policy of a firm is held constant, its dividend payout policy has no consequences for shareholder wealth. Higher dividend payouts lead to lower retained earnings and capital gains, and vice versa, leaving total wealth of the shareholders unchanged. Contrary to this prediction, however, corporations follow extremely deliberate dividend payout strategies (Lintner (1956)). This evidence raises a puzzle: How do firms choose their dividend policies?

In the United States and other countries, the puzzle is even deeper since many shareholders are taxed more heavily on their dividend receipts than on capital gains. The actual magnitude of this tax burden is debated (see Poterba and Summers (1985) and Allen and Michaely (1997)), but taxes generally make it even harder to explain dividend policies of firms.

Economists have proposed a number of explanations of the dividend puzzle. Of these, particularly popular is the idea that firms can signal future profitability by paying dividends (see Bhattacharya (1979), John and Wil-

1

liams (1985), Miller and Rock (1985), and Ambarish, John, and Williams (1987)). Empirically, this theory had considerable initial success, since firms that initiate (or raise) dividends experience share price increases, and the converse is true for firms that eliminate (or cut) dividends (Aharony and Swary (1980), Asquith and Mullins (1983)). Recent results are more mixed, since current dividend changes do not help predict firms' future earnings growth (DeAngelo, DeAngelo, and Skinner (1996) and Benartzi, Michaely, and Thaler (1997)).

Another idea, which has received only limited attention until recently (e.g., Easterbrook (1984), Jensen (1986), Fluck (1998, 1999), Hart and Moore (1974), Myers (1998), Gomes (2000), and Zwiebel (1996)), is that dividend policies address agency problems between corporate insiders and outside shareholders. According to these theories, unless profits are paid out to shareholders, they may be diverted by the insiders for personal use or committed to unprofitable projects that provide private benefits for the insiders. As a consequence, outside shareholders have a preference for dividends over retained earnings. Theories differ on how outside shareholders actually get firms to disgorge cash. The key point, however, is that failure to disgorge cash leads to its diversion or waste, which is detrimental to outside shareholders' interest.

The agency approach moves away from the assumptions of the Modigliani–Miller theorem by recognizing two points. First, the investment policy of the firm cannot be taken as independent of its dividend policy, and, in particular, paying out dividends may reduce the inefficiency of marginal investments. Second, and more subtly, the allocation of all the profits of the firm to shareholders on a pro rata basis cannot be taken for granted, and in particular the insiders may get preferential treatment through asset diversion, transfer prices, and theft—even holding the investment policy constant. Insofar as dividends are paid on a pro rata basis, they benefit outside shareholders relative to the alternative of expropriation of retained earnings.

In this paper, we attempt to identify some of the basic elements of the agency approach to dividends, to understand its key implications, and to evaluate them on a cross section of more than 4,000 firms from 33 countries around the world. The reason for looking around the world is that the severity of agency problems to which minority shareholders are exposed differs greatly across countries, in part because legal protection of these shareholders varies (La Porta, Lopez-de-Silanes, Shleifer, and Vishny (1997, 1998), henceforth referred to as LLSV). Empirically, we find that dividend policies vary across legal regimes in ways consistent with a particular version of the agency theory of dividends. Specifically, firms in common law countries, where investor protection is typically better, make higher dividend payouts than firms in civil law countries do. Moreover, in common but not civil law countries, high growth firms make lower dividend payouts than low growth firms. These results support the version of the agency theory in which investors in good legal protection countries use their legal powers to extract dividends from firms, especially when reinvestment opportunities are poor.

Section I of the paper summarizes some of the theoretical arguments. Section II describes the data. Section III presents our empirical findings. Section IV concludes.

## I. Theoretical Issues

### A. *Agency Problems and Legal Regimes*

Conflicts of interest between corporate insiders, such as managers and controlling shareholders, on the one hand, and outside investors, such as minority shareholders, on the other hand, are central to the analysis of the modern corporation (Berle and Means (1932), Jensen and Meckling (1976)). The insiders who control corporate assets can use these assets for a range of purposes that are detrimental to the interests of the outside investors. Most simply, they can divert corporate assets to themselves, through outright theft, dilution of outside investors through share issues to the insiders, excessive salaries, asset sales to themselves or other corporations they control at favorable prices, or transfer pricing with other entities they control (see Shleifer and Vishny (1997) for a discussion). Alternatively, insiders can use corporate assets to pursue investment strategies that yield them personal benefits of control, such as growth or diversification, without benefiting outside investors (e.g., Baumol (1959), Jensen (1986)).

What is meant by *insiders* varies from country to country. In the United States, the U.K., Canada, and Australia, where ownership in large corporations is relatively dispersed, most large corporations are to a significant extent controlled by their managers. In most other countries, large firms typically have shareholders that own a significant fraction of equity, such as the founding families (La Porta, Lopez-de-Silanes, and Shleifer (1999)). The controlling shareholders can effectively determine the decisions of the managers (indeed, managers typically come from the controlling family), and hence the problem of *managerial* control per se is not as severe as it is in the rich common law countries. On the other hand, the controlling shareholders can implement policies that benefit themselves at the expense of minority shareholders. Regardless of the identity of the insiders, the victims of insider control are minority shareholders. It is these minority shareholders who would typically have a taste for dividends.

One of the principal remedies to agency problems is the law. Corporate and other law gives outside investors, including shareholders, certain powers to protect their investment against expropriation by insiders. These powers in the case of shareholders range from the right to receive the same per share dividends as the insiders, to the right to vote on important corporate matters, including the election of directors, to the right to sue the company for damages. The very fact that this legal protection exists probably explains why becoming a minority shareholder is a viable investment strategy, as opposed to just being an outright giveaway of money to strangers who are under few if any obligations to give it back.

As pointed out by LLSV (1998), the extent of legal protection of outside investors differs enormously across countries. Legal protection consists of both the content of the laws and the quality of their enforcement. Some countries, including most notably the wealthy common law countries such as the United States and the U.K., provide effective protection of minority shareholders so that the outright expropriation of corporate assets by the insiders is rare. Agency problems manifest themselves primarily through non-value-maximizing investment choices. In many other countries, the condition of outside investors is a good deal more precarious, but even there some protection does exist. LLSV (1998) show in particular that common law countries appear to have the best legal protection of minority shareholders, whereas civil law countries, and most conspicuously the French civil law countries, have the weakest protection.

The quality of investor protection, viewed as a proxy for lower agency costs, has been shown to matter for a number of important issues in corporate finance. For example, corporate ownership is more concentrated in countries with inferior shareholder protection (LLSV (1998), La Porta, Lopez-de-Silanes, and Shleifer (1999)). The valuation and breadth of capital markets is greater in countries with better investor protection (LLSV (1997), Demirguc-Kunt and Maksimovic (1998)). Finally, there is some evidence that good investor protection contributes to the efficiency of resource allocation and to economic growth more generally (Levine and Zervos (1998), Rajan and Zingales (1995)). This paper continues this research by examining the dividend puzzle using shareholder protection as a proxy for agency problems.

## B. Agency and Dividends: Two Views

### B.1. The Role of Dividends in an Agency Context

In a world of significant agency problems between corporate insiders and outsiders, dividends can play a useful role. By paying dividends, insiders return corporate earnings to investors and hence are no longer capable of using these earnings to benefit themselves. Dividends (a bird in the hand) are better than retained earnings (a bird in the bush) because the latter might never materialize as future dividends (can fly away). Additionally, the payment of dividends exposes companies to the possible need to come to the capital markets in the future to raise external funds, and hence gives outside investors an opportunity to exercise some control over the insiders at that time (Easterbrook (1984)).

Unfortunately, there are no fully satisfactory theoretical agency models of dividends that derive dividend policies as part of some broad optimal contract between investors and corporate insiders, which allows for a range of feasible financing instruments. Instead, different models, such as Fluck (1998, 1999), Myers (1998), and Gomes (2000), capture different aspects of the problem. Moreover, the existing agency models do not fully deal with the issues of choice between debt and equity in addressing agency problems, the

choice between dividends and share repurchases, and the relationship between dividends and new share issues. We attempt to distill from the available literature the basic mechanisms of how dividends could be used to deal with agency problems. In particular, we distinguish between two very different agency "models" of dividends. The predictions of these models that we test are necessarily limited by the fact that we do not look at all the financing and payout choices simultaneously.

Perhaps most importantly in this regard, we do not examine share repurchases, which have been commonly taken as an alternative to paying dividends. We note, however, that share repurchases are most common precisely in the countries where firms pay high dividends, such as the United States and the U.K. For example, between June 1997 and June 1998 there were 1,537 share repurchases in the world recorded by the Securities Data Corporation, of which 1,100 occurred in the United States. By market value, the United States accounted for 72 percent of world share repurchases during this period, and the United States, the U.K., Canada, and Australia combined accounted for 83 percent. In some civil law countries, share repurchases are even illegal or heavily taxed (*The Economist*, August 15, 1998).[1] If share repurchases are complementary to dividends, rather than a substitute for them, our evidence only underestimates the difference in total cash payouts to shareholders between civil and common law countries.

## B.2. Dividends as an Outcome of Legal Protection of Shareholders

Under the first view, dividends are an outcome of an effective system of legal protection of shareholders. Under an effective system, minority shareholders use their legal powers to force companies to disgorge cash, thus precluding insiders from using too high a fraction of company earnings to benefit themselves.[2] Shareholders might do so by voting for directors who offer better dividend policies, by selling shares to potential hostile raiders who then gain control over non–dividend paying companies, or by suing companies that spend too lavishly on activities that benefit only the insiders. Moreover, good investor protection makes asset diversion legally riskier and more expensive for the insiders, thereby raising the relative attraction of dividends for them. The greater the rights of the minority shareholders, the more cash they extract from the company, other things equal.

It is important to recognize that this argument does not rely on minority shareholders having specific rights to dividends per se, but rather on their having the more general rights of voting for directors and protesting wealth

---

[1] It could be argued that the discouragement of share repurchases is a form of shareholder protection since, unlike dividends, share repurchases can be discriminatory. This argument is less plausible in light of the fact that most share repurchases in the United States and the U.K. are open market, and, moreover, appear to supplement rather than substitute for dividends.

[2] Even under an effective system, residual agency problems must remain, for if they are totally resolved, we are back to the world of Modigliani and Miller with no reason for dividends.

expropriation. A good example from the United States is Kirk Kerkorian forcing Chrysler Corporation to disgorge its cash by paying dividends in 1995 to 1996. As a large shareholder in Chrysler, Kerkorian had no specific rights to dividends, but used the voting mechanism to put his associates on the board and then force the board to sharply raise dividends. Another good example is Velcro Industries, the producer of the famous "touch fastener" incorporated on the island of Curaçao in the Netherlands Antilles, "where shareholders have no right of dissent" (*Forbes*, October 15, 1990). Two-thirds of the shares of Velcro Industries are controlled by the Cripps family that runs Velcro (*Forbes*, May 23, 1994). In 1988, despite having a large cash reserve, the company suspended dividends "for the foreseeable future" (*Forbes*, October 3, 1988), delisted itself from the Montreal Stock Exchange, and aggressively wrote down assets to slash earnings, evidently to "buy out Velcro minority holders cheap" (*Forbes*, May 23, 1994). The share price dived and, in 1990, with dividends remaining at zero, the Crippses offered to repurchase minority shares at slightly above the market price. Minority shareholders sued in New York and "when a New York judge ruled that the United States was the proper jurisdiction, secretive Sir Humphrey Cripps decided to call off his offer rather than go under the light of U.S. court of law" (*Forbes*, May 23, 1994). The company subsequently resumed its dividend payments. This case illustrates that, in a high protection country like the United States, in contrast to a low protection country like the Netherlands, shareholders are able to extract dividends from companies by virtue of their ability to resist oppression rather than having any specific dividend rights per se.

In a cross section of countries with different quality of shareholder protection, the implication that better protection is associated with higher dividend payouts is testable. There is one further implication of this theory. Consider a country with good shareholder protection, and compare two companies in that country: one with good investment opportunities and growth prospects, and another with poor opportunities. Shareholders who feel protected would accept low dividend payouts, and high reinvestment rates, from a company with good opportunities because they know that when this company's investments pay off, they could extract high dividends. In contrast, a mature company with poor investment opportunities would not be allowed to invest unprofitably. As a consequence, with good shareholder protection, high growth companies should have significantly lower dividend payouts than low growth companies. In contrast, if shareholder protection is poor, we would not necessarily expect such a relationship between payouts and growth since shareholders may try to get what they can—which may not be much—immediately. This also is a testable implication.[3] The implications of the outcome agency model of dividends are illustrated in Figure 1.

---

[3] Ambarish et al. (1987) derive the negative relationship between growth and payouts in a dividend signaling model. They do not focus on how this relationship would vary depending on how well shareholders are protected. In principle, this extension is possible.

**Figure 1. Outcome model of dividends.**

## B.3. *Dividends as a Substitute for Legal Protection of Shareholders*

In an alternative agency view, dividends are a substitute for legal protection.[4] This view relies crucially on the need for firms to come to the external capital markets for funds, at least occasionally. To be able to raise external funds on attractive terms, a firm must establish a reputation for moderation in expropriating shareholders. One way to establish such a reputation is by paying dividends, which reduces what is left for expropriation. For this mechanism to work, the firm must never want to "cash in" its reputation by stopping dividends and expropriating shareholders entirely. The firm would never want to cash in if, for example, there is enough uncertainty about its future cash flows that the option of going back to the capital market is always valuable (Bulow and Rogoff (1989)).

A reputation for good treatment of shareholders is worth the most in countries with weak legal protection of minority shareholders, who have little else to rely on. As a consequence, the need for dividends to establish a reputation is the greatest in such countries. In countries with stronger shareholder protection, in contrast, the need for a reputational mechanism is weaker, and hence so is the need to pay dividends. This view implies that, other things equal, dividend payout ratios should be higher in countries with weak legal protection of shareholders than in those with strong protection.[5]

---

[4] The closest informal discussion to the substitute model is Easterbrook (1984). Formally, the model that comes the closest to taking this point of view is Gomes (2000). However, the recent drafts of his paper have moved away from focusing on dividends, and hence our discussion should not be interpreted as a description of Gomes's model.

[5] Dewenter and Warther (1998) argue that there is less need to signal future earnings with dividends in Japan than in the United States. This may be because Japanese firms have better ways of information transmission to the relevant investors than do U.S. firms, or because Japanese managers are more insulated from investor pressure (Kang and Stulz (1996)). Dewenter and Warther find that share price reactions to dividend changes are smaller in Japan than in the United States. This finding may be consistent with either of the two agency models of dividends.

**Figure 2. Substitute model of dividends.**

Additionally, in this view, firms with better growth prospects also have a stronger incentive to establish a reputation since they have a greater potential need for external finance, other things equal. As a result, firms with better growth prospects might choose higher dividend payout ratios than firms with poor growth prospects. However, firms with good growth prospects also have a better current use of funds than firms with poor growth prospects. The relationship between growth prospects and dividend payout ratios is therefore ambiguous. Figure 2 illustrates the implications of this substitute agency model of dividends.

### B.4. Summary of Predictions of Agency Models

We refer to the two alternative agency models of dividends as "the outcome model" and "the substitute model." The outcome model predicts that dividend payout ratios are higher in countries with good shareholder protection, other things equal. The substitute model predicts the opposite. The outcome model further predicts that, in countries with good shareholder protection, companies with better investment opportunities should have lower dividend payout ratios. The substitute model does not make this prediction. In fact, it makes a weak prediction that, in countries with poor shareholder protection, firms with better investment opportunities might pay out more to maintain reputations.

### C. Tax Issues

Economists are divided on the effects of taxes on the valuation of dividends (Poterba and Summers (1985)). The so-called traditional view holds that heavy taxation of dividends at both the corporate and personal levels—at least in the United States—is a strong deterrent to paying out dividends rather than retaining the earnings. There are two important objections to this view. One objection, raised by Miller and Scholes (1978), states that investors have

access to a variety of dividend tax avoidance strategies that allow them to effectively escape dividend taxes. This objection does not closely correspond to what investors actually do (Feenberg (1981)). Another objection, the so-called new view of dividends and taxes (e.g., King (1977), Auerbach (1979)), holds that cash has to be paid out as dividends sooner or later, and therefore paying it earlier in the form of current dividends imposes no greater a tax burden on shareholders than does the delay. According to this theory, taxes do not deter dividend payments. Harris, Hubbard, and Kemsley (1997) support this new view. In our empirical work, we include a measure of the tax disadvantage of dividends based on Poterba and Summers (1984, 1985) to assess the effect of taxes on dividend policies. Appendix A summarizes in detail our treatment of the tax effects of dividends, and also presents the data on taxes that we use in the empirical work.

## II. Data

Our sample is based on the March 1996 edition of the WorldScope Database, which presents information on the (typically) largest listed firms in 46 countries. There are 13,698 firms in the original database. Since accounting data are often reported with a delay, our analysis uses data through 1994. Table I, Panel A summarizes the construction of the sample. From the original universe, we eliminate firms trading in socialist countries and in Luxembourg; firms listed in countries with mandatory dividend policies (i.e., legal requirements that a certain fraction of net income is paid out as dividends); financial firms; firms completely or partially owned by the government (as best we can identify them); firms without consolidated balance sheets in 1989, 1994, or both; firms with negative net income or negative cash flow in 1994; firms with missing dividend data in 1994 or missing sales, net income, or cash flow data in 1994 or 1989; firms whose dividends exceed sales; and finally, three firms that do not appear to be publicly traded. This leaves us with the basic sample of 4,103 firms from 33 countries for which we can compute dividend payout ratios in 1994 and sales growth rates from 1989 to 1994. Panel B shows how we get from 46 to 33 countries.

We note in particular the exclusion of countries with mandatory dividend rules, namely Brazil, Chile, Colombia, Greece, and Venezuela.[6] Some of these countries have weak legal protection of minority shareholders. The fact that, in such environments, regulators choose to force companies to pay dividends is in itself some evidence in favor of the importance of agency considerations, since the most plausible reason for a mandatory dividend policy is to assure outside investors that they would not be expropriated entirely, and thus to encourage participation in the equity markets by such investors (LLSV (1998)). In general, firms in mandatory dividend countries have higher pay-

---

[6] There also appears to be a minimum dividend requirement in Germany, although it can be waived at the discretion of management. Because this requirement is so weak, we do not count Germany as a mandatory-dividend country. Excluding it would only strengthen our results.

**Table I**
**Construction of the Sample**

| | Panel A: Firms in the Sample |
|---:|---|
| 13,698 | WorldScope Sample (3/96 version) |
| −56 | Firms listed in stock exchanges of former socialist countries |
| −12 | Firms listed in Luxembourg's stock exchange |
| −323 | Firms listed in stock exchanges of countries with mandatory dividend policies |
| −2,836 | Financial firms (primary and/or secondary SIC between 6,000 and 6,999) |
| −335 | State-owned enterprises (direct and/or indirect government ownership) |
| −1,296 | Unconsolidated balance sheets in 1989, 1994, or both |
| −3,878 | Missing sales in 1989 and/or dividends, cash flows, net income or sales in 1994 |
| −832 | Negative net income before extraordinary items in 1994 |
| −11 | Negative cash flow in 1994 |
| −13 | Dividends > Sales |
| −3 | Not publicly traded (i.e., cooperatives and privately owned firms) |
| 4,103 | Basic sample |
| | Panel B: Countries in the Sample |
| 46 | Countries in WorldScope |
| −3 | Socialist, former socialist countries (China, Poland, Hungary) |
| −1 | Luxembourg |
| −5 | Mandatory dividend countries (Brazil, Chile, Colombia, Greece, Venezuela) |
| −4 | Countries that do not meet data requirements (Israel, Pakistan, Peru, Sri Lanka) |
| 33 | Countries in the sample |

outs than firms in countries without such rules, but they nevertheless appear, in the data, to have lower payouts than required by the law. A possible reason for this is that the accounting earnings reported to the authorities for the purposes of compliance with mandatory dividend rules are lower than the earnings reported to the shareholders which we use in our analysis.

Table II summarizes the construction of the variables. We use two rough proxies for protection of minority shareholders. The first is a dummy equal to one if a country's company law or commercial code is of civil origin, and zero for common law origin. Because we have data on few countries, we do not distinguish between French, German, and Scandinavian civil law origins in this paper, as in LLSV (1997, 1998). In general, civil law countries have weaker legal protection of minority shareholders than do common law countries. The second measure of investor protection, the low investor protection dummy, is equal to one if the index of antidirector rights is below the sample median. The index of antidirector rights comes from LLSV (1998), and reflects such aspects of minority rights as the ease of voting for directors, the possibility of electing directors through a cumulative voting mechanism, the existence of a grievance mechanism for oppressed minority shareholders, such as a class action lawsuit, the percentage of votes needed to call an extraordinary shareholder meeting, and the existence of preemptive rights.

Since we are dealing with accounting data in countries with different accounting standards, we compute several measures of the dividend payout ratio. The numerator in these ratios is the total cash dividend paid to common and preferred shareholders. The denominators are cash flow, earnings, and sales. The dividend-to-cash-flow ratio has a natural economic interpretation since it is the ratio of cash distributed to cash generated in a period. The dividend-to-earnings ratio is the most commonly used measure of dividend payouts. The two ratios have several problems, however. First, both of them may depend on a country's accounting conventions, and hence may not be exactly comparable across countries. Second, these ratios have the potential problem of being easily manipulated by accounting tricks. Third, and perhaps most important, diversion of resources may occur before earnings or cash flows are reported, in which case these two ratios overestimate the share of true earnings that is paid out as dividends. Fortunately, if diversion is greater in countries with poor shareholder protection, this problem biases the results toward finding higher payouts in these countries than is really the case. Our results of lower measured payouts in countries with poor shareholder protection reported below would thus be even stronger if true earnings and cash flows were higher than reported. Still, as an additional guard against these problems, we also present the dividends-to-sales ratio, since sales are less dependent on accounting conventions, are harder to manipulate or smooth through accounting practices, and are less subject to theft. Sales should be viewed just as a deflator; the economic interpretation of this ratio is not transparent.

The trickiest measurement problem we face is how to capture investment opportunities across firms in a way that is consistent across countries. Our principal measure of such opportunities is the past growth in sales of each firm, which has the advantage of being roughly independent of accounting practices, but has the disadvantage of relying on the past as a proxy for the future. For each firm, we compute its annual real sales growth rate over the five-year period from 1989 to 1994. In Section III, we discuss other measures of investment opportunities.

For our dividend payout ratios and the sales growth rate, we also compute industry-adjusted measures. For each company in a given industry, we make this adjustment relative to the worldwide rather than countrywide measure for that industry (i.e., we take out worldwide industry effects rather than country-industry effects). Consider the computation of the industry-adjusted growth in sales, for example. We first find for each industry in each country the median real growth rate of sales in that industry in that country. We then take the median of country medians, thus obtaining the worldwide median growth in real sales in the industry. Our measure of industry-adjusted growth in sales for a company is the difference between that company's sales growth and the world median sales growth in its industry. The idea is that different industries might be at different stages of maturity and growth that determine their dividend policies.

**Table II**
## The Variables

This table describes the variables collected for the 33 countries included in our study. The first column gives the name of the variable, the second column describes the variable and provides the sources for the variables.

| Variable | Description |
|---|---|
| Common law | Equals one if the origin of the Company Law or Commercial Code of the country is the English Common Law and zero otherwise. *Source:* LLSV (1998). |
| Civil law | Equals one if the Company Law or Commercial Code of the country originates in Roman Law and zero otherwise. *Source:* LLSV (1998). |
| Low protection | Equals one if the index of antidirectors rights is smaller or equal to three (the sample median) and zero otherwise. The index of antidirectors rights is formed by adding one when: (1) the country allows shareholders to mail their proxy vote; (2) shareholders are not required to deposit their shares prior to the General Shareholders' Meeting; (3) cumulative voting or proportional representation of minorities on the board of directors is allowed; (4) an oppressed minorities mechanism is in place; (5) the minimum percentage of share capital that entitles a shareholder to call for an Extraordinary Shareholders' Meeting is less than or equal to 10 percent (the sample median); (6) or when shareholders have preemptive rights that can only be waved by a shareholders meeting. The range for the index is from zero to six. *Source*: LLSV (1998). |
| High protection | Equals one if the index of antidirectors rights (defined above) is greater than three (the sample median) and zero otherwise. *Source*: LLSV (1998). |
| Dividend-to-cash-flow | Dividends as a percentage of cash flow in fiscal year 1994. Dividends are defined as total cash dividends paid to common and preferred shareholders. Cash flow is measured as total funds from operations net of non-cash items from discontinued operations. *Source*: WorldScope Database. |
| IA_dividend-to-cash-flow | Industry-adjusted dividend-to-cash-flow ratio for a firm. To calculate IA_dividend-to-cash-flow, we first find for each industry in each country the median of the dividend-to-cash-flow ratio (C_D/CF). Then for each industry in the sample we define the *world median* as the median of C_D/CF across countries. Finally, we calculate IA_dividend-to-cash-flow as the difference between the firm's dividend-to-cash-flow and the *world* median dividend-to-cash-flow for the firm's industry. We rely on a firm's primary SIC to define the following seven broad industries: (1) agriculture; (2) mining; (3) construction; (4) light manufacturing; (5) heavy manufacturing; (6) communications and transportation; and (7) services. *Source*: WorldScope Database. |
| Dividend-to-earnings | Dividends as a percentage of earnings in fiscal year 1994. Dividends are defined as total cash dividends paid to common and preferred shareholders. Earnings are measured after taxes and interest but before extraordinary items. *Source*: WorldScope Database. |

| | |
|---|---|
| IA_dividend-to-earnings | Industry-adjusted dividend-to-earnings ratio for a firm. To calculate IA_dividend-to-earnings, we first find for each industry in each country the median of the dividend-to-earnings ratio (C_D/E). Then for each industry in the sample we define the world median as the median of C_D/E across countries. Finally, we calculate IA_dividend-to-earnings as the difference between the firm's dividend-to-earnings and the world median dividend-to-earnings for the firm's industry. We rely on a firm's primary SIC to define the following seven broad industries: (1) agriculture; (2) mining; (3) construction; (4) light manufacturing; (5) heavy manufacturing; (6) communications and transportation; and (7) services. *Source*: WorldScope Database. |
| Dividend-to-sales | Dividends as a percentage of sales in fiscal year 1994. Dividends are defined as total cash dividends paid to common and preferred shareholders. Sales are net sales. *Source*: WorldScope Database. |
| IA_dividend-to-sales | Industry-adjusted dividend-to-sales ratio for a firm. To calculate IA_dividend-to-sales, we first find for each industry in each country the median of the dividend-to-sales ratio (C_D/S). Then for each industry in the sample we define the world median as the median of C_D/S across countries. Finally, we calculate IA_dividend-to-sales as the difference between the firm's dividend-to-sales and the world median dividend-to-sales for the firm's industry. We rely on a firm's primary SIC to define the following seven broad industries: (1) agriculture; (2) mining; (3) construction; (4) light manufacturing; (5) heavy manufacturing; (6) communications and transportation; and (7) services. *Source*: WorldScope Database. |
| GS | Average annual percentage growth in real (net) sales over the period 1989–1994. Before computing GS, we translate net sales in U.S. dollars into real terms by using the U.S. GNP deflator. *Source*: WorldScope Database and *International Financial Statistics* (1996). |
| GS_decile | Rank decile for GS. Firms are ranked by legal origin into 10 equal-size groups. Ranges from 1 to 10 in ascending order of GS. |
| IA_GS | Average annual industry-adjusted growth in (net) sales over the period 1989–1994. To calculate IA_GS, we first find for each industry in each country the median of the GS (C_GS). Then for each industry in the sample we define the world median as the median of C_GS across countries. Finally, we calculate IA_GS as the difference between the firm's GS and the world median GS for the firm's industry. We rely on a firm's primary SIC to define the following seven broad industries: (1) agriculture; (2) mining; (3) construction; (4) light manufacturing; (5) heavy manufacturing; (6) communications and transportation; and (7) services. *Source*: WorldScope Database. |
| IA_GS_decile | Rank decile for IA_GS. It ranges from 1 to 10. |
| Dividends tax advantage | The ratio of the value, to an outside investor, of US$1 distributed as dividend income to the value of US$1 received in the form of capital gains when kept inside the firm as retained earnings. The computation of this ratio is described in Appendix A. *Sources*: Ernst and Young's *Worldwide Corporate Tax Guide and Directory* (1994), Price Waterhouse's *Individual Taxes: A Worldwide Summary* (1995), and OECD's *Taxing Profits in a Global Economy: Domestic and International Issues* (1991). |

<div align="center">

**Table III**

**The Data**

</div>

Panel A classifies countries by legal origin and presents medians by country. Definitions for each of the variables can be found in Table II. Panel B reports tests of medians for civil versus common legal origin.

| Country | N | Low Protection | Div/CF (%) | Div/Earn (%) | Div/Sales (%) | GS (Annual) | Div Tax Adv. |
|---|---|---|---|---|---|---|---|
| | | | Panel A: Medians | | | | |
| Argentina | 3 | 0 | 12.65 | 27.36 | 4.32 | 14.32 | 1.00 |
| Austria | 9 | 1 | 5.85 | 24.83 | 0.77 | 13.31 | 0.78 |
| Belgium | 33 | 1 | 11.77 | 39.38 | 1.09 | 3.78 | 0.74 |
| Denmark | 75 | 1 | 6.55 | 17.27 | 0.71 | 4.32 | 0.67 |
| Finland | 39 | 1 | 8.08 | 21.27 | 0.77 | −2.14 | 1.07 |
| France | 246 | 1 | 9.46 | 23.55 | 0.63 | 4.54 | 0.64 |
| Germany | 146 | 1 | 12.70 | 42.86 | 0.83 | 5.88 | 0.86 |
| Indonesia | 1 | 1 | 8.72 | 25.11 | 0.77 | 32.62 | 0.76 |
| Italy | 58 | 1 | 9.74 | 21.83 | 0.92 | −1.38 | 0.77 |
| Japan | 149 | 0 | 13.03 | 52.88 | 0.72 | 6.19 | 0.70 |
| South Korea | 2 | 1 | 7.33 | 18.49 | 0.66 | 5.29 | 0.79 |
| Mexico | 14 | 1 | 19.47 | 46.44 | 3.59 | 8.02 | 1.00 |
| Netherlands | 96 | 1 | 11.29 | 30.02 | 0.74 | 4.13 | 0.40 |
| Norway | 50 | 0 | 10.74 | 23.91 | 0.98 | 4.43 | 1.08 |
| Philippines | 4 | 1 | 6.72 | 10.47 | 2.45 | −7.29 | 1.05 |
| Portugal | 17 | 1 | 0.64 | 38.01 | 0.64 | 8.20 | 0.98 |
| Spain | 33 | 0 | 15.77 | 30.45 | 1.04 | 1.32 | 0.72 |
| Sweden | 81 | 1 | 5.59 | 18.33 | 0.78 | −0.63 | 1.03 |
| Switzerland | 70 | 1 | 10.38 | 25.30 | 0.98 | 3.73 | 0.56 |
| Taiwan | 3 | 1 | 48.97 | 68.89 | 11.54 | 1.62 | 0.60 |
| Turkey | 6 | 1 | 8.61 | 22.64 | 2.08 | 0.16 | 0.90 |
| **Civil Law Median** | **33** | **1** | **9.74** | **25.11** | **0.83** | **4.32** | **0.78** |
| Australia | 103 | 0 | 22.83 | 42.82 | 2.22 | 2.21 | 0.90 |
| Canada | 236 | 0 | 8.00 | 19.78 | 0.78 | −0.62 | 0.89 |
| Hong Kong | 40 | 0 | 35.43 | 45.93 | 7.51 | 7.94 | 1.00 |
| India | 1 | 0 | 25.69 | 49.34 | 1.55 | −0.09 | 0.59 |
| Ireland | 16 | 0 | 17.39 | 27.28 | 0.96 | 9.96 | 0.77 |
| Malaysia | 41 | 0 | 15.29 | 37.93 | 3.12 | 16.31 | 0.68 |
| New Zealand | 17 | 0 | 19.16 | 35.60 | 2.26 | 3.11 | 1.00 |
| Singapore | 27 | 0 | 22.28 | 41.04 | 2.14 | 11.02 | 0.96 |
| South Africa | 90 | 0 | 16.16 | 35.62 | 1.90 | 3.47 | 0.85 |
| Thailand | 10 | 1 | 32.83 | 52.56 | 3.35 | 17.73 | 0.90 |
| United Kingdom | 799 | 0 | 16.67 | 36.91 | 1.89 | 2.44 | 0.83 |
| United States | 1,588 | 0 | 11.38 | 22.11 | 0.95 | 3.15 | 0.58 |
| **Common Law Median** | **40** | **0** | **18.28** | **37.42** | **2.02** | **3.31** | **0.87** |
| **Sample Median** | **39** | **1** | **11.77** | **30.02** | **0.98** | **4.13** | **0.83** |
| | | | Panel B: Test of Medians (*z*-statistic) | | | | |
| Civil vs Common Law | | 3.97* | −3.29* | −1.72*** | −2.36** | −0.34 | −0.09 |

*, **, and *** indicate significance at the 1, 5, and 10 percent levels, respectively.

Table III summarizes the data by presenting the number of observations we have for each country as well as country medians of several variables. Of the firms in our sample, a little over one-quarter (1,135) are from civil law countries and a little over three-quarters (2,968) are from common law countries. More than half of the firms in the sample come from the United States

and the United Kingdom. Both of these countries have a large number of listed firms; WorldScope coverage and the quality of data are also better for richer countries. India, for example, has 5,398 listed firms in 1995, but only one of them makes it into the sample.

The second column of Table III illustrates the finding of our earlier work, namely that common law countries on average have stronger shareholder protection, as illustrated by the median of the low shareholder protection dummy, than do civil law countries. The *z*-statistic on the difference in the median civil law and common law shareholder protection is 3.97.

The next three columns present country medians of our three dividend payout ratios. The median of country median dividend-to-earnings ratios (the most common payout metric used in the United States) is about 30 percent, confirming that a substantial share of earnings is paid out as dividends.[7] Paying dividends is indeed what large firms just about everywhere do, and there is a dividend puzzle to be explained. Table III also reveals that, for all measures, common law countries have higher payouts than civil law countries, and for two out of three the difference is statistically significant at the 5 percent level. We discuss this result in more detail below.

The sixth column shows that the median of country medians real growth rate of sales in the sample is 4.13 percent. At the median of country medians, firms in civil law countries grow one percent faster than firms in common law countries.

A final point in Table III is that, in most countries, the difference between the tax treatment of dividends and retained earnings is small. The United States, with its significant tax advantage of retained earnings, is relatively extreme.[8]

## III. Results

### A. Simple Statistics

We present the results in three steps. First, in Tables IV and V, we present some basic statistics from our sample of firms that bear on the hypotheses described in Section I. In computing these statistics, we weigh all the countries equally, so the United States and the U.K., where most firms in the sample are located, do not receive any extra weight. Second, in Tables VI and VII, we present the regressions on a cross section of companies that control for tax and industry effects. In these regressions, countries that have more companies automatically receive more weight. These two ways of presenting the data are thus complementary, since one can argue for both empirical strategies. Finally, we discuss the robustness of our results to several alternative measurement and specification strategies.

---

[7] Note that, in the calculation of this measure, the United States and the U.K. do not receive any more weight than any other country.

[8] In the computation of tax rates, we combine federal and local taxes. For example, for the United States we add federal (28 percent) and New York State (7.75 percent) capital gains tax rates.

<div align="center">

**Table IV**

**Dividends by Legal Origin and Growth Opportunities**

</div>

This table classifies firms based on both the legal origin of the country in which they are incorporated and on their growth in sales (GS) relative to the world median growth in sales. Countries are required to have at least five valid observations (firms) with growth in sales below the world median and five observations with growth in sales above the world median. The number of countries in the resulting sample is 24 (14 civil law and 10 common law countries). To compute the world median growth in sales we calculate the median growth in sales for each country and then we take medians again but now over the 24 resulting country observations. For each classification, the table reports the median value of the country medians for the following three ratios: (1) dividend-to-cash-flow in Panel A; (2) dividend-to-earnings in Panel B; and (3) dividend-to-sales in Panel C. Finally, Panel D reports $Z$-statistics for tests of difference in medians.

| Legal Origin | All | "Growth" GS> World Median GS | "Mature" GS< World Median GS |
|---|---|---|---|
| *Panel A: Dividend-to-cash-flow* | | | |
| Civil law | 10.56 | 10.89 | 9.20 |
| Common law | 17.03 | 15.17 | 22.87 |
| *Panel B: Dividend-to-earnings* | | | |
| Civil law | 27.66 | 30.35 | 21.27 |
| Common law | 36.27 | 27.95 | 40.88 |
| *Panel C: Dividend-to-sales* | | | |
| Civil law | 0.80 | 0.89 | 0.77 |
| Common law | 2.02 | 1.77 | 2.91 |
| *Panel D: Z-statistic for Differences in Medians* | | | |
| | Div/CF | Div/Earn | Div/Sales |
| Civil vs Common law | −2.81* | −0.76 | −2.75* |
| Civil law: Mature vs growth | −0.92 | −0.87 | −0.92 |
| Common law: Mature vs growth | 2.34** | 2.42** | 1.74*** |

*, **, and *** indicate significance at the 1, 5, and 10 percent levels, respectively.

In Tables IV and V, we present medians of country medians (MOMs) of dividend payout ratios for various groups of firms, and in particular distinguish between rapidly and slowly growing firms. To have reasonably robust statistics, we use a narrower sample in these tables than we do in Table III. Specifically, we only consider countries where we have at least five firms with sales growth above the world median sales growth of 4.1 percent, and five firms with sales growth below the world median. This restriction leaves us with 24 countries, and eliminates countries with very few firms from the analysis. In the regressions, we go back to the broader sample.[9]

In Table IV, we examine whether firms in civil and common law countries have different dividend payout policies. To begin, we compute the MOM for the three dividend payout ratios for the civil and common law families sep-

---

[9] We have also computed the medians without the restriction on the number of firms with high and low growth rates in each country. The results are very similar.

<div align="center">

**Table V**

**Dividends by Legal Protection and Growth Opportunities**

</div>

This table classifies firms based both on the level of investor protection of the country in which they are incorporated (low or high protection) and on their growth in sales (GS) relative to the world median growth in sales. Countries included are required to have at least five valid observations (firms) with growth in sales below the world median and five observations with growth in sales above the world median. The number of countries in the resulting sample is 24 (11 with low protection equal to one). To compute the world median growth in sales we calculate the median growth in sales for each country and then we take medians again but now over the 24 resulting country-observations. For each classification, the table reports the median value of the country-medians for the following three ratios: (1) dividend-to-cash-flow in Panel A; (2) dividend-to-earnings in Panel B; and (3) dividend-to-sales in Panel C. Finally, Panel D reports $Z$-statistics for tests of difference in medians.

| Investor Protection | All | "Growth" GS>World Median GS | "Mature" GS<World Median GS |
|---|---|---|---|
| Panel A: Dividend-to-cash-flow | | | |
| Low protection | 9.74 | 10.86 | 8.74 |
| High protection | 16.16 | 14.51 | 18.93 |
| Panel B: Dividend-to-earnings | | | |
| Low protection | 25.30 | 31.31 | 21.24 |
| High protection | 35.62 | 29.05 | 39.69 |
| Panel C: Dividend-to-sales | | | |
| Low protection | 0.78 | 0.88 | 0.76 |
| High protection | 1.89 | 1.53 | 2.24 |
| Panel D: $Z$-statistic for Differences in Medians | | | |
| | Div/CF | Div/Earn | Div/Sales |
| Low vs high protection | −2.87* | −1.13 | −2.40** |
| Low protection: Mature vs growth | −1.15 | −1.08 | −1.54 |
| High protection: Mature vs growth | 2.38** | 2.23** | 1.67*** |

*, **, and *** indicate significance at the 1,5, and 10 percent levels, respectively.

arately (the same measures, for a broader sample, are presented in Table III). The results of this calculation are presented in the first column of Table IV. For all three ratios, common law countries have a higher dividend payout ratio than civil law countries do. The MOM dividend-to-cash-flow ratio is 17 percent for common law countries, and only 10.6 percent for civil law countries. The MOM dividend-to-earnings ratios are 36.3 percent for common law countries, and 27.7 percent for civil law countries. The MOM dividend-to-sales ratio is two percent for common law countries and 0.8 percent for civil law countries. For all three variables, these estimates are very close to those for the broader sample in Table III. Panel D of Table IV shows that, for two out of the three measures of dividend payouts, the difference between the common law MOM payout and the civil law MOM payout is statistically significant.

The results in the first column of Table IV are central to this paper. Recall from Table III that common law countries generally have stronger minority shareholder protection than civil law countries. The fact that common law

countries also have higher dividend payouts supports the outcome agency model of dividends, according to which better shareholder protection leads to higher dividend payouts. In contrast, the result is inconsistent with the basic prediction of the substitute agency model of dividends. More generally, the fact that dividend payouts are so different in environments with different shareholder protection suggests that agency considerations are likely to be central to the explanation of why firms pay dividends.

The additional results in Table IV address the relationship between dividend payout rates and sales growth rates across legal regimes. For each country with enough observations (see above), we separately compute the median payout ratio for firms with above and firms with below the world median sales growth rate. Within each origin, we then compute the MOM payout across countries for rapidly and slowly growing firms separately. The results are presented in the last two columns of Table IV, and again are consistent across all three measures of dividend payouts. In common law countries, payout ratios are strictly higher for slowly growing firms than for rapidly growing firms. In the common law family, the MOM dividend-to-cash-flow ratio is 15.2 percent for rapidly growing firms and 22.9 percent for slowly growing firms; the MOM dividend-to-earnings ratio is 28 percent for rapidly growing firms and 41 percent for slowly growing firms; and the MOM dividend-to-sales ratio is 1.8 percent for rapidly growing firms and 2.9 percent for slowly growing firms. These differences between mature and growth firms in common law countries are statistically significant (see Panel D). These results are consistent with the predictions of the outcome agency model, according to which well-protected minority shareholders are willing to delay dividends in firms with good growth prospects.

In the civil law family, in contrast, rapidly growing firms appear, if anything, to pay higher dividends. In this family, the MOM dividend-to-cash-flow ratio is 10.9 percent for rapidly growing firms and 9.2 percent for slowly growing firms; the MOM dividend-to-earnings ratio is 30.3 percent for rapidly and 21.3 percent for slowly growing firms; and finally the MOM dividend-to-sales ratio is 0.9 percent for rapidly and 0.8 percent for slowly growing firms. The positive association between dividend payouts and growth rates in civil law countries is consistent with the dividends as substitutes theory applying to these countries. However, as Panel D shows, these payout differences between mature and growth firms in civil law countries are not statistically significant, and hence we should not read too much into this finding.

Table V presents calculations similar to those in Table IV, except that now countries are sorted by whether the low shareholder protection dummy is equal to zero or one. As in Table IV, we use the narrow sample of countries. The results are similar to those in Table IV, and we summarize them only briefly. First, on all measures of dividend payouts, countries with better shareholder protection have higher dividend payout ratios than do countries with worse protection. Second, again on all measures of dividend payouts, within countries with good shareholder protection, high growth firms have

lower dividend payouts than low growth firms. The differences are statistically significant at the 5 percent level in two cases, and at the 10 percent level in the third. Finally, on all measures of dividend payouts, within countries with low shareholder protection, high growth firms have higher dividend payouts than low growth firms. These differences are not statistically significant, however.

The preliminary results are consistent with the outcome agency model. However, the findings may be driven by some heterogeneity of countries correlated with legal origin or investor protection. Accordingly, we next move to a regression analysis that attempts to control for the differences in tax regimes and in industrial composition in different countries.

## B. Regressions

Table VI presents the results of regressions across 4,103 firms in 33 countries around the world. We use the broader sample described in Table III. We employ a random effects specification that explicitly accounts for the cross-correlation between error terms for firms in the same country. We control for the tax advantage of dividends, which is specific to each country, but not for industry effects until Table VII. We report results for all three measures of the dividend payout ratio. We use dummies to proxy for the quality of legal protection of investors. For each payout variable, we present one regression that distinguishes between common and civil law countries, and one regression that distinguishes between low and high shareholder protection countries, and one that includes both the origin and the protection dummies. As a measure of investment opportunities in the regressions, we use the decile rank of the past average annual sales growth rate for each firm, GS_decile. In this calculation, the deciles of growth rates are defined separately for companies in common and civil law families. Using deciles gives us a less widely spread variable, and defining deciles separately for the two families ensures that we have enough high growth firms in civil law countries. We also include an interaction between GS_decile and the legal origin or the low investor protection dummy.

The tax variable enters with the positive sign in all specifications, but is only statistically significant in the dividend-to-sales ratio regressions. The interpretation of this result is highly ambiguous. The positive coefficients can be interpreted as some support for the traditional view, under which taxes discourage the payment of dividends. The insignificance of these coefficients, however, may be interpreted as evidence in favor of the "new view," under which tax payments are already capitalized in the value of the firm and therefore do not influence dividend policy. Last, the evidence may mean that our computations do not adequately address the nuances of each country's tax treatment of dividends.

Consider first the regressions that use only one measure of investor rights at a time. The civil law dummy enters with a negative and significant coefficient at the 1 percent level in regressions using all three measures of

## Table VI
## Regression Results for Raw Data

Regressions with country random effects for the cross section of 33 countries around the world. The dependent variables are the 1994 values of the following three ratios: (1) dividend-to-cash-flow; (2) dividend-to-earnings; and (3) dividend-to-sales. The independent variables are: (1) civil law, a dummy variable that equals one if the legal origin of the Company Law or Commercial Code of the country in which the firm is incorporated is Roman Law and zero otherwise; (2) low protection, a dummy variable that equals one if the Index of Antidirectors rights (described in Table II) of the country in which the firm is incorporated is equal or smaller than three (the sample median) and zero otherwise; (3) GS, the firm's average annual percentage growth in sales over the period 1989 to 1994; (4) the interaction between GS and civil law origin; (5) the interaction between GS and Low Protection; and (6) tax advantage of retained earnings (described in Table II). Standard errors are shown in parenthesis.

| | | | Dependent Variables | | | | | |
| Constant | Civil law | Low protection | GS_decile | GS_decile*Civil | GS_decile* Low protection | Div tax advantage | N | $\chi^2$ |
|---|---|---|---|---|---|---|---|---|
| | | | Panel A: Dividend-to-cash-flow as Dependent Variable | | | | | |
| 22.3730* | −13.2591* | | −0.8457* | 0.9022* | | 3.2262 | 4,103 | 137.79* |
| (3.5145) | (1.6602) | | (0.0832) | (0.1608) | | (3.9635) | | |
| 20.2817* | | −10.8156* | −0.8133* | | 0.8554* | 3.5303 | 4,103 | 109.52* |
| (5.0539) | | (2.1943) | (0.0813) | | (0.1695) | (5.7621) | | |
| 22.6043* | −13.2883* | −0.1404 | −0.8502* | 0.8948* | 0.1112 | 3.1591 | 4,103 | 134.29* |
| (3.8440) | (3.0909) | (3.1446) | (0.0832) | (0.3504) | (0.3676) | (4.3237) | | |
| | | | Panel B: Dividend-to-earnings as Dependent Variable | | | | | |
| 44.1156* | −16.4633* | | −2.1354* | 2.3925* | | 9.5905 | 4,102 | 119.38* |
| (8.4626) | (3.9817) | | (0.1974) | (0.3816) | | (9.5425) | | |
| 44.6786* | | −18.1518* | −2.0613* | | 2.4253* | 8.9278 | 4,102 | 118.13* |
| (8.4796) | | (4.0195) | (0.1927) | | (0.4007) | (9.7170) | | |
| 44.9493* | −8.1284 | −10.1918 | −2.1431* | 1.5135*** | 1.0124 | 8.9453 | 4,102 | 121.61* |
| (8.9882) | (7.2780) | (7.3403) | (0.1974) | (0.8308) | (0.8717) | (10.1111) | | |
| | | | Panel C: Dividend-to-sales as Dependent Variable | | | | | |
| 1.8963* | −2.0821* | | −0.0859* | 0.0962* | | 1.8157* | 4,103 | 157.59* |
| (0.4005) | (0.2185) | | (0.0142) | (0.0273) | | (0.4556) | | |
| 1.4299 | | −1.6413* | −0.0884* | | 0.0926* | 2.1266** | 4,103 | 61.76* |
| (0.7812) | | (0.3461) | (0.0139) | | (0.0288) | (0.8911) | | |
| 1.8907* | −2.3471* | 0.2979 | −0.0865* | 0.1189** | −0.0248 | 1.8457* | 4,103 | 153.45* |
| (0.4146) | (0.4378) | (0.4476) | (0.0142) | (0.0595) | (0.0623) | (0.4708) | | |

*, **, and *** indicate significance at the 1, 5, and 10 percent levels, respectively.

**Table VII**

## Regression Results for Industry-Adjusted Data

Regressions with country random effects for the cross-section of thirty three countries around the world. The dependent variables are the 1994 values of the following three ratios: (1) industry-adjusted-dividend-to-cash-flow; (2) industry-adjusted-dividend-to-earnings; and (3) industry-adjusted-dividend-to-sales. The independent variables are: (1) civil law, a dummy variable that equals one if the origin of the Company Law or Commercial Code of the country in which the firm is incorporated is Roman Law and zero otherwise; (2) low protection, a dummy variable that equals one if the Index of Antidirectors rights (described in Table II) of the country in which the firm is incorporated is equal or smaller than three (the sample median) and zero otherwise; (3) IA_GS, the firm's annual average percentage industry-adjusted growth in sales over the period 1989–1994; (4) the interaction between IA_GS and civil law origin; (5) the interaction between IA_GS and low protection; and (6) tax advantage of retained earnings (calculated as indicated in Table II). We require at least five observations in each country/industry and report only on the industries that have the required number of observations in at least three countries. Standard errors are shown in parenthesis.

| | | | Dependent Variables | | | | | |
|---|---|---|---|---|---|---|---|---|
| Constant | Civil law | Low protection | IA_GS_decile | IA_GS_decile* Civil | IA_GS_decile* Low protection | Div tax advantage | N | $\chi^2$ |
| *Panel A: Industry-adjusted-dividend-to-cash-flow as Dependent Variable* | | | | | | | | |
| 10.0288* | −12.7246* | | −0.8730* | 0.8869* | | 3.7703 | 4,077 | 141.26* |
| (3.6114) | (1.6883) | | (0.0826) | (0.1598) | | (4.0724) | | |
| 8.1130*** | | −10.5460* | −0.8343* | | 0.8295* | 4.0120 | 4,077 | 117.12* |
| (4.8398) | | (2.1150) | (0.0806) | | (0.1688) | (5.5198) | | |
| 10.2997* | −12.4283* | −0.5290 | −0.8758* | 0.8946* | −0.0102** | 3.6333 | 4,077 | 138.58* |
| (3.9021) | (3.0636) | (3.0865) | (0.0827) | (0.3413) | (0.3592) | (4.3890) | | |
| *Panel B: Industry-adjusted-dividend-to-earnings as Dependent Variable* | | | | | | | | |
| 14.2369 | −15.9368* | | −2.2892* | 2.4323* | | 7.7032 | 4,076 | 134.59* |
| (9.7285) | (4.4173) | | (0.1980) | (0.3835) | | (10.9624) | | |
| 14.4038 | | −16.8178* | −2.1865* | | 2.3746* | 7.2467 | 4,076 | 129.53* |
| (9.7308) | | (4.2450) | (0.1932) | | (0.4041) | (11.1096) | | |
| 14.9785 | −10.8174 | −6.3909 | −2.2938* | 1.9432** | 0.5648 | 7.1076 | 4,076 | 135.48* |
| (10.2583) | (7.7188) | (7.7540) | (0.1981) | (0.8191) | (0.8622) | (11.5312) | | |
| *Panel C: Industry-adjusted-dividend-to-sales as Dependent Variable* | | | | | | | | |
| 1.1042* | −2.1146* | | −0.1087* | 0.1244* | | 1.4371* | 4,077 | 147.28* |
| (0.4248) | (0.2238) | | (0.0139) | (0.0268) | | (0.4819) | | |
| 0.6490 | | −1.6415* | −0.1076* | | 0.1165* | 1.6823*** | 4,077 | 77.82* |
| (0.8000) | | (0.3507) | (0.0135) | | (0.0283) | (0.9124) | | |
| 1.1116* | −2.5125* | 0.4255 | −0.1106* | 0.1557* | −0.0332 | 1.4931* | 4,077 | 138.28* |
| (0.4715) | (0.4439) | (0.4520) | (0.0139) | (0.0571) | (0.0600) | (0.5329) | | |

*, **, and *** indicate significance at the 1, 5, and 10 percent levels, respectively.

dividend payouts.[10] Using the dividend to cash flow ratio, for example, common law countries have a 13.3 percentage point higher payout, other things equal. The coefficient on GS_decile is negative and also significant at the 1 percent level, and implies that, for common law countries, moving from the bottom to the top decile of sales growth rate is associated with a 7.6 percentage point lower dividend to cash flow ratio. That is, in common law countries, higher growth firms pay moderately lower dividends. At the same time, the coefficient on the interaction between GS_decile and the civil law dummy is highly statistically significant and of roughly the same magnitude as that on GS_decile in all three regressions. This implies that, other things equal, there is no relationship between sales growth and dividend payouts in civil law countries. The results using the civil law dummy, like the medians in Table IV, are consistent with the outcome agency model of dividends.[11]

Similar results obtain using the low shareholder protection dummy. The coefficient on that dummy is negative and significant at the 1 percent level using all measures of payout.[12] The coefficient on GS_decile as before is negative and significant, implying that, in countries with good shareholder protection, faster growing firms pay lower dividends. The coefficient on the interaction between GS_decile and the low shareholder protection dummy is positive and of about the same magnitude, indicating that the relationship between growth and payouts does not hold in countries with poor shareholder protection. These results also suggest that dividends are an outcome of pressure on the insiders to pay out profits.

When both the civil law dummy and the poor shareholder protection dummy are included in the regression, in two out of three cases the former remains significant, while the latter does not. (In the third case, both variables lose significance.) Although it is best not to put too much weight on this result given that the two variables are correlated, one view is that our measure of shareholder protection does not perfectly capture some of the differences between the legal regimes. For example, as argued in LLSV (1998), the quality of law enforcement—which surely matters for shareholder power—is also better in common law than in civil law countries. The other results do not change appreciably when both dummies are included at the same time.

In Table VII, we use industry-adjusted growth in sales and industry-adjusted dividends to control for industry effects, and otherwise estimate the same equations as in Table VI (the details of the adjustment are described in Table II). The industry adjustment does not change the thrust of

[10] The civil law dummy is also highly significant when included in the regression on its own, without the growth in sales variables.
[11] These results also survive the inclusion of a measure of the quality of accounting standards, described in LLSV (1998), available for 31 countries in the sample (not Ireland and Indonesia).
[12] The poor shareholder protection dummy is also highly significant when included in the regression on its own, without the growth in sales variables.

our results. Countries from the common law family, as well as countries with good shareholder protection, pay higher industry-adjusted dividends, and, moreover, in these countries, faster growing firms pay lower dividends, other things equal.

## C. Robustness

In this subsection, we briefly describe the results of some of the robustness checks of our findings. One question is whether the regression results are shaped by firms from the United States and the U.K., which are the majority of the sample. Of course, the results in Tables IV and V weigh all countries equally, but one might want to know more about firm-level data. Accordingly, Figures 3 and 4 present the plots of dividend payouts against sales growth for each of the 11 common law and 20 civil law countries respectively.[13] Figure 3 shows that there is a negative relationship between growth in sales and dividend-to-earnings ratios in every one of the 11 common law countries. Figure 4 shows that this relationship is negative for 11 of the 20 civil law countries, and positive for nine of the 20. If we plot the ratio of dividends to cash flow against sales growth, the relationship is again negative for all 11 common law countries, and for 11 out of 20 civil law countries. Finally, if we plot the ratio of dividends to sales against sales growth, the relationship is negative for 10 of the 11 common law countries, and for 10 of the 20 civil law countries. In summary, while the results for different countries hold with different levels of statistical significance, they consistently show that more rapidly growing firms pay lower dividends in common law, but not in civil law countries.[14]

A further concern about our results is that we might have selected a particular point in time during national (or international) business cycles that makes our results special. To address this concern, we reestimate all regressions using 1992, 1993, and 1994 dividend variables, and look at three-year rather than five-year past sales growth rates (thus, for example, we have related measures of 1992 dividends to 1989 to 1991 sales growth rates). Our results hold using these alternative points in time for measuring dividend payouts and investment opportunities.

A related point deals with the inherent crudeness in measuring investment opportunities in terms of the past growth rate in sales. We have chosen to use the past growth rate in sales to avoid the incompatibility of accounting variables across countries. To check robustness, we have also reestimated our results using growth rates of assets, fixed assets, cash flow, and earnings, as well as industry $Q$, as measures of investment opportuni-

---

[13] We do not have enough observations to run a regression for India and Indonesia.

[14] Very similar results obtain if we divide three countries by high versus low antidirector rights.

**Figure 3. Dividends-to-earnings ratios for common law countries**. Scatter plots are shown of dividend-to-earnings ratios (div/earn) against growth-in-sales (GS) for 11 common law countries (India does not have a plot because it has only one observation). To avoid outliers, we cap the maximum dividend-to-earnings ratio at the common law 95th percentile.
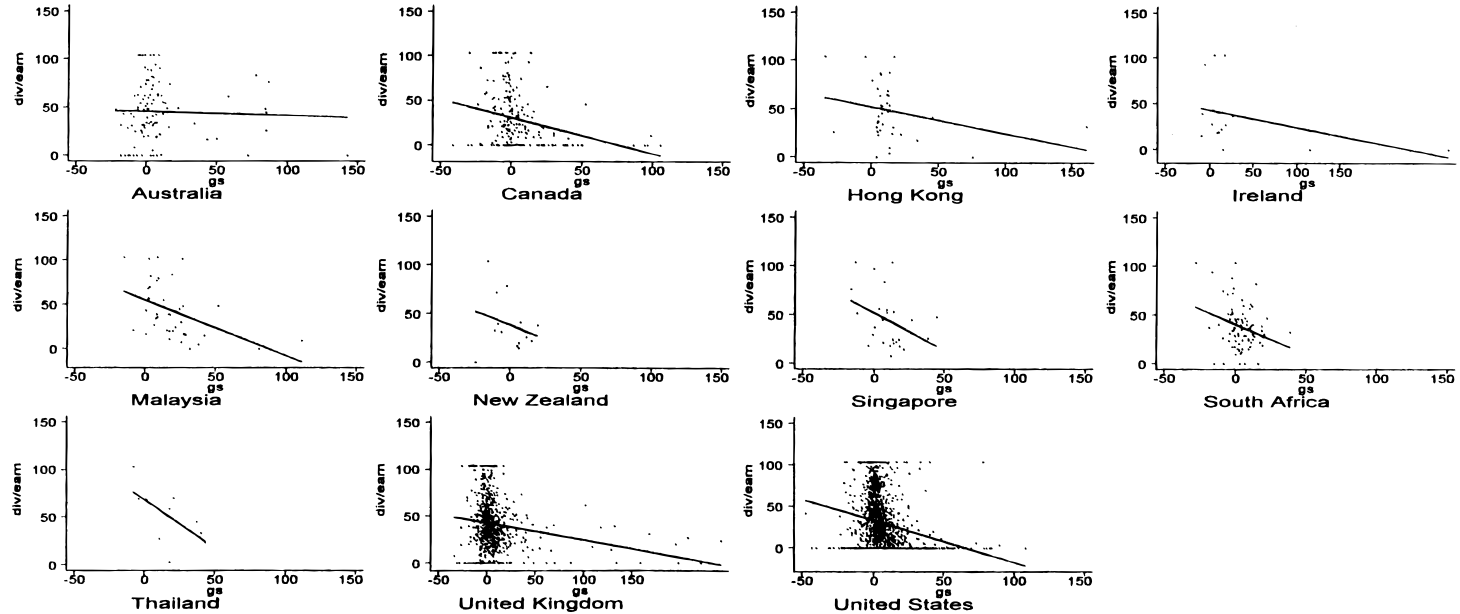
**Figure 4. Dividends-to-earnings ratios for civil law countries**. Scatter plots are shown of dividend-to-earnings ratios (div/earn) against growth-in-sales (GS) for 20 civil law countries (Indonesia does not have a plot because it has only one observation). To avoid outliers, we cap the maximum dividend-to-earnings ratio at the civil law 95th percentile.

ties. The results generally confirm the reported findings in both sign and significance, although the relationship between industry $Q$ and dividends is insignificant.

One possible alternative interpretation of our results is that our measures of investor protection simply reflect the degree of capital market development. It is possible that firms in developed capital markets are happy to pay out their earnings because they can always raise more external funds, whereas firms in undeveloped capital markets would hold on to the hard-to-get cash. This view would explain our finding that, on average, dividend payouts are higher in countries with good investor protection, which also happen to be countries with developed capital markets.

This alternative view has its own problems, however. To begin, the degree of capital market development is to a significant extent endogenous, and indeed in part determined by legal origin and the quality of investor protection (LLSV (1997)). Moreover, this view does not explain our findings on the relationship between investment opportunities and payouts. If anything, this view would imply that firms in poorly developed capital markets should exhibit extreme sensitivity of payouts to growth opportunities, and really try to hoard cash when they have good investments. In contrast, firms in developed markets should be willing to pay dividends regardless of investment opportunities since they can count on raising external funds. Contrary to these predictions, our data show that the negative relationship between investment opportunities and payouts is stronger in countries with good investor protection and hence more developed capital markets.

As a final point, we briefly address a possibly important objection to our analysis, which states that perhaps the evidence of lower payouts in civil law (or poor shareholder protection) countries simply reflects greater reliance on debt finance in those countries. First, as an empirical matter, we use the ratios of dividends to cash flow and to earnings, so the denominators already take out interest payments. Even if firms in civil law countries rely on debt to a greater extent, they should not necessarily pay out less of their net-of-interest income. Second, it is not generally the case that firms in civil law countries rely more on debt finance. Indeed, many of these countries, particularly French civil law countries, have poor legal protection of both shareholders and creditors, and hence have both smaller debt and smaller equity markets (LLSV (1997)). The idea that countries with poorly developed stock markets necessarily, or even on average, have better developed lending mechanisms is simply a myth. Last, we actually test the validity of this objection by including a country-specific measure of debt finance from LLSV (1997), namely the ratio of aggregate private debt to GNP, in the regressions in Tables VI and VII. The coefficients on the debt variable are positive, though generally insignificant, while the magnitudes and the statistical significance of shareholder protection coefficients remain largely unaffected. This finding is inconsistent with the argument that poor shareholder protection is associated with lower dividend payouts because of substitution of financing into debt.

## IV. Conclusion

This paper uses a sample of firms from 33 countries around the world to shed light on dividend policies of large corporations. We take advantage of different legal protection of minority shareholders across these countries to compare dividend policies of companies whose minority shareholders face different risks of expropriation of their wealth by corporate insiders. We use this cross-sectional variation to examine the agency approach to dividend policy.

We distinguish two alternative agency models of dividends. In the first model, dividends are an outcome of effective legal protection of shareholders, which enables minority shareholders to extract dividend payments from corporate insiders. In the second, dividends are a substitute for effective legal protection, which enables firms in unprotective legal environments to establish reputations for good treatment of investors through dividend policies.

Our data suggest that the agency approach is highly relevant to an understanding of corporate dividend policies around the world. More precisely, we find consistent support for the outcome agency model of dividends. Firms operating in countries with better protection of minority shareholders pay higher dividends. Moreover, in these countries, fast growth firms pay lower dividends than slow growth firms, consistent with the idea that legally protected shareholders are willing to wait for their dividends when investment opportunities are good. On the other hand, poorly protected shareholders seem to take whatever dividends they can get, regardless of investment opportunities. This apparent misallocation of investment is presumably part of the agency cost of poor legal protection.

In our analysis, we find no conclusive evidence on the effect of taxes on dividend policies. Nor can we use our data to assess the relevance of dividend signaling. In fact, our results are consistent with the idea that, on the margin, dividend policies of firms may convey information to some investors. Despite the possible relevance of alternative theories, firms appear to pay out cash to investors because the opportunities to steal or misinvest it are in part limited by law, and because minority shareholders have enough power to extract it. In this respect, the quality of legal protection of investors is as important for dividend policies as it is for other key corporate decisions.

## Appendix A

Table A.I presents the raw data used to calculate the tax preference of dividends for each country. We use the tax rates faced by local residents who acquire minority stakes in publicly traded securities and hold their investments long enough to qualify for long-term capital gains tax rates. Furthermore, we assume that the effective tax rate on capital gains is equivalent to one-fourth of the nominal rate (Poterba (1987)). Finally, we combine federal and local taxes whenever possible. In order to compute the tax parameter, it is helpful to use the criteria proposed by King (1977) and group the tax systems of the countries in our sample in three broad categories:

**Table A.I**
**Construction of the Tax Advantage of Dividends**

| Country | (A) Corporate Tax Undistributed Profits | (B) Corporate Tax Distributed Profits | (C) Personal Tax Capital Gains | (D) Personal Tax Dividends | (E) Imputation Rate | (G) Value of $1 in Dividends $(1 - B + E) * (1 - D)$ | (H) Value of $1 in Capital Gains $(1 - A) * (1 - C/4)$ | Dividend Tax Preference (G/H) |
|---|---|---|---|---|---|---|---|---|
| Argentina | 0.30 | 0.30 | 0.00 | 0.00 | 0.00 | 0.70 | 0.70 | 1.00 |
| Austria | 0.34 | 0.34 | 0.00 | 0.22 | 0.00 | 0.51 | 0.66 | 0.78 |
| Belgium[1] | 0.40 | 0.40 | 0.00 | 0.26 | 0.00 | 0.44 | 0.60 | 0.74 |
| Denmark | 0.34 | 0.34 | 0.40 | 0.40 | 0.00 | 0.40 | 0.59 | 0.67 |
| Finland | 0.25 | 0.25 | 0.25 | 0.25 | 0.25 | 0.75 | 0.70 | 1.07 |
| France[2] | 0.33 | 0.33 | 0.19 | 0.60 | 0.33 | 0.40 | 0.63 | 0.63 |
| Germany[3] | 0.54 | 0.41 | 0.00 | 0.53 | 0.25 | 0.39 | 0.46 | 0.86 |
| Indonesia[4] | 0.35 | 0.35 | 0.30 | 0.30 | 0.00 | 0.46 | 0.60 | 0.76 |
| Italy[5] | 0.52 | 0.52 | 0.00 | 0.51 | 0.27 | 0.37 | 0.48 | 0.77 |
| Japan[6] | 0.52 | 0.52 | 0.26 | 0.35 | 0.00 | 0.31 | 0.45 | 0.70 |
| S. Korea[7] | 0.34 | 0.34 | 0.00 | 0.22 | 0.00 | 0.52 | 0.66 | 0.79 |
| Mexico[8] | 0.41 | 0.41 | 0.00 | 0.00 | 0.00 | 0.59 | 0.59 | 1.00 |
| Netherlands | 0.35 | 0.35 | 0.00 | 0.60 | 0.00 | 0.26 | 0.65 | 0.40 |
| Norway | 0.28 | 0.28 | 0.28 | 0.28 | 0.28 | 0.72 | 0.67 | 1.08 |
| Philippines | 0.35 | 0.35 | 0.20 | 0.00 | 0.00 | 0.65 | 0.62 | 1.05 |
| Portugal[9] | 0.40 | 0.40 | 0.10 | 0.30 | 0.22 | 0.57 | 0.59 | 0.97 |
| Spain | 0.35 | 0.35 | 0.56 | 0.56 | 0.26 | 0.40 | 0.56 | 0.72 |
| Sweden | 0.28 | 0.28 | 0.13 | 0.00 | 0.00 | 0.72 | 0.70 | 1.03 |
| Switzerland[10] | 0.34 | 0.34 | 0.00 | 0.44 | 0.00 | 0.37 | 0.66 | 0.56 |
| Taiwan | 0.25 | 0.25 | 0.00 | 0.40 | 0.00 | 0.45 | 0.75 | 0.60 |
| Turkey[11] | 0.27 | 0.27 | 0.00 | 0.10 | 0.00 | 0.66 | 0.73 | 0.90 |
| **Civil Law Mean** | **0.36** | **0.35** | **0.13** | **0.30** | **0.09** | **0.52** | **0.62** | **0.81** |
| Australia | 0.33 | 0.33 | 0.47 | 0.47 | 0.33 | 0.53 | 0.59 | 0.90 |
| Canada[12] | 0.44 | 0.44 | 0.40 | 0.36 | 0.14 | 0.45 | 0.51 | 0.89 |
| Hong Kong | 0.18 | 0.18 | 0.00 | 0.00 | 0.00 | 0.83 | 0.83 | 1.00 |
| India[13] | 0.52 | 0.52 | 0.22 | 0.45 | 0.00 | 0.27 | 0.46 | 0.58 |
| Ireland | 0.40 | 0.40 | 0.40 | 0.48 | 0.20 | 0.42 | 0.54 | 0.77 |
| Malaysia[14] | 0.30 | 0.30 | 0.00 | 0.32 | 0.00 | 0.48 | 0.70 | 0.68 |
| New Zealand | 0.33 | 0.33 | 0.00 | 0.33 | 0.33 | 0.67 | 0.67 | 1.00 |
| Singapore | 0.27 | 0.27 | 0.00 | 0.30 | 0.27 | 0.70 | 0.73 | 0.96 |
| South Africa[15] | 0.40 | 0.49 | 0.00 | 0.00 | 0.00 | 0.51 | 0.60 | 0.85 |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Thailand | 0.30 | 0.30 | 0.00 | 0.37 | 0.30 | 0.63 | 0.70 | 0.90 |
| United Kingdom | 0.33 | 0.33 | 0.40 | 0.40 | 0.17 | 0.50 | 0.60 | 0.83 |
| United States[16] | 0.42 | 0.42 | 0.36 | 0.47 | 0.00 | 0.31 | 0.53 | 0.58 |
| **Common Law Mean** | **0.35** | **0.36** | **0.19** | **0.33** | **0.14** | **0.53** | **0.62** | **0.85** |

*Notes*:

[1]Corporate tax rates in Belgium include a three percent crisis contribution surtax. The corporate rate is 39 percent.

[2]Dividends in France are grossed up by 50 percent for tax purposes and the individual can claim credit for up to 50 percent of the cash amount of the dividend. Personal taxes on dividends include 56.8 percent of income tax and 3.4 percent of social contribution. Personal taxes on capital gains tax are calculated as the sum of the 16 percent basic rate and 3.4 percent of social contribution.

[3]Dividends in Germany are grossed up by 3/7 for tax purposes and the individual can claim credit for up to 3/7 of the cash amount of the dividend. Municipal tax rates on corporate income range from 13 percent to 19 percent (16 percent average used here) and are deductible.

[4]Personal capital gains in Indonesia are taxed as ordinary income (30 percent).

[5]Dividends in Italy are grossed up by 56.25 percent for tax purposes and the individual can claim credit for up to 56.25 percent of the cash amount of the dividend. Corporate taxes are the sum of 36 percent corporate income tax (IRPEG) and 16.2 percent local income tax (ILOR).

[6]Corporate income tax in Japan is calculated as the sum of three terms: (1) 37.5 percent corporate income tax; (2) 20.7 percent surcharge (Tokyo metropolitan area); and (3) 13.2 percent enterprise tax (deductible).

[7]Corporate taxes in Korea include a 7.5 percent resident tax surcharge on top of the 32 percent corporate tax rate.

[8]Corporate taxes in Mexico include a 10 percent mandatory employee-profit-sharing contribution (deductible) in addition to the 34 percent corporate tax rate.

[9]Corporate taxes in Portugal include a 10 percent municipal surcharge (derrama) in addition to the 36 percent basic rate. The tax rate of 30 percent on dividends distributed by SA corporations includes five percent inheritance tax.

[10]Combined cantonal and communal corporate tax rates range from 21.7 percent to 46.65 percent in Switzerland. We took the middle point for corporate taxes. We used average combined local and federal for personal dividend tax rates.

[11]Corporate taxes in Turkey include a seven percent surtax in addition to the basic corporate tax rate (25 percent).

[12]Dividends in Canada are grossed up by 25 percent for tax purposes and the individual can claim credit for up to 25.0 percent of the cash amount of the dividend. The 14 percent imputation rate is based on the highest combined federal and provincial marginal tax rates for individuals in Ontario (35.92 percent). Corporate taxes include both a three percent surtax as well as a 15.5 percent provincial tax (Ontario) in addition to the basic rate (28 percent). Personal capital gains and dividend taxes are the maximum combined federal and provincial marginal tax rates for Ontario residents.

[13]Indian corporate taxes are based on a 45 percent basic rate and a 15 percent surcharge. Similarly, the personal dividend and capital gains tax of 20 percent is augmented by a 12 percent surcharge.

[14]Capital gains taxes are not adjusted for a sales tax of 0.25 percent on each trade.

[15]Corporate taxes on distributed profits in South Africa include a 15 percent surtax (secondary tax on companies or STC) on dividends declared or paid after March 17, 1993 on the top of the 40 percent corporate tax rate.

[16]The U.S. corporate tax rate includes a 6.5 percent (average) local tax rate on top of the 35 percent federal tax rate. The individual capital gains and dividend taxes those applicable to residents of the state of New York (7.875 percent).

*Sources*: *Worldwide Corporate Tax Guide and Directory*, Ernst and Young, 1994.

   *Worldwide Personal Tax Guide*, Ernst and Young, 1994.

   *Corporate Taxes, A Worldwide Summary*, Price Waterhouse, 1995.

   *Individual Taxes*, *A Worldwide Summary*, Price Waterhouse, 1995.

   *Taxing Profits in a Global Economy*, *Domestic and International Issues*, OECD, 1991.

Whenever Ernst and Young and Price Waterhouse differed, we rely on the source that presents more details. We use the OECD source only for Switzerland.

1. The Classical System: Personal and corporate taxation are independent of each other and shareholders receive no compensation for taxes paid at the corporate level. Specifically, the company pays a flat rate of corporate tax on profits (i.e., distributed and undistributed income are taxed at the same rate) and shareholders pay income tax on dividend receipts. Accordingly, the value to an investor of one dollar in earnings distributed in the form of dividends is equal to $(1 - \tau_{\text{corp}}) * (1 - \tau_{\text{div}})$, where $\tau_{\text{corp}}$ is the corporate tax rate on income and $\tau_{\text{div}}$ is the personal tax rate on dividend receipts. Similarly, the value to an investor of one dollar in earnings retained inside the firm is given by $(1 - \tau_{\text{corp}}) * (1 - \tau_{\text{cap}})$, where $\tau_{\text{cap}}$ is the effective personal tax rate on capital gains. Therefore, the dividend tax preference parameter (defined as the ratio of the value earnings distributed as dividends versus earnings retained inside the firm) is given by $(1 - \tau_{\text{div}})/(1 - \tau_{\text{cap}})$.

2. *The Two-Rate System*: The corporate tax rate on earnings distributed as dividends is lower than on retained earnings to mitigate the tax advantage of retained earnings in the classical system. Accordingly, the value to an investor of one dollar in earnings distributed in the form of dividends is equal to $(1 - \tau_{\text{dist}}) * (1 - \tau_{\text{div}})$, where $\tau_{\text{dist}}$ is the corporate tax rate on distributed income. Similarly, the value to an investor of one dollar in earnings retained inside the firm is given by $(1 - \tau_{\text{ret}}) * (1 - \tau_{\text{cap}})$, where $\tau_{\text{ret}}$ is the corporate tax rate on retained earnings. Thus, the dividend tax preference parameter is given by $(1 - \tau_{\text{dist}}) * (1 - \tau_{\text{div}})/((1 - \tau_{\text{ret}}) * (1 - \tau_{\text{cap}}))$. In practice, the pure two-rate system is implemented rarely in our sample of countries. In fact, only two countries in our sample have different tax rates for retained earnings and dividends: Germany and South Africa. However, in South Africa the taxes on dividends are higher than on retained earnings contrary to the motivation behind the *two-rate system*, ($\tau_{\text{div}} = 49$ percent versus $\tau_{\text{ret}} = 40$ percent). Interestingly, dividends in Germany are not only taxed at a lower corporate rate but shareholders are allowed to credit taxes paid by corporations on distributions to offset personal taxes in the same way as in the *imputation system*.

3. *The Imputation System*: Shareholders receive credit for taxes paid by the company on earnings distributed as dividends. These credits may be used to offset shareholder's tax liability. Part of the corporate tax liability on distributed profits is "imputed" to shareholders and regarded as a prepayment of their personal income tax. In the most frequent version of the *imputation system*, dividends are regarded as having borne personal tax at the "imputation" rate $\tau_{\text{imp}}$ and shareholders are liable only for the difference between their marginal tax rates on personal income and the imputation rate (i.e., they pay taxes on dividend receipts at the rate $\tau_{\text{div}} - \tau_{\text{imp}}$). Accordingly, the value to an investor of one dollar in earnings distributed in the form of dividends is equal to $(1 - \tau_{\text{corp}} + \tau_{\text{imp}}) * (1 - \tau_{\text{div}})$. Hence, the dividend tax preference parameter is given by $(1 - \tau_{\text{dist}} + \tau_{\text{imp}}) * (1 - \tau_{\text{div}})/((1 - \tau_{\text{ret}}) * (1 - \tau_{\text{cap}}))$.

Less frequently, the operation of the system is defined in terms of a tax credit rate $\tau_{\text{cred}}$ and not an imputation rate. In countries that rely on tax credits, shareholders are liable for the difference between the personal taxes owed on dividends-cum-tax-credit received and the tax credit (i.e., they pay taxes on dividend receipts at the rate $(1 + \tau_{\text{cred}}) * \tau_{\text{div}} - \tau_{\text{cred.}}$). In such cases, we re-express $\tau_{\text{cred}}$ in terms of its associated $\tau_{\text{imp}}$ and use the formula for the *imputation system*.

## Appendix B

Summary statistics of the data in the paper are presented in Table B.I. The variables are defined in Table II.

**Table B.I**
**Summary Statistics**

| Variable | Observations | Mean | Median | Standard Deviation | Minimum | Maximum |
|---|---|---|---|---|---|---|
| Civil law | 4,103 | 0.2766 | 0 | 0.4474 | 0 | 1 |
| Low protection | 4,103 | 0.2218 | 0 | 0.4155 | 0 | 1 |
| Dividend-to-cash-flow | 4,103 | 15.1209 | 13.0677 | 13.4504 | 0 | 47.7362 |
| IA_dividend-to-cash-flow | 4,077 | 3.2143 | 1.1791 | 13.2962 | −14.7019 | 34.3852 |
| Dividend-to-earnings | 4,103 | 35.2640 | 29.8006 | 31.9907 | 0 | 134.3036 |
| IA_dividend-to-earnings | 4,076 | 3.5192 | −2.0261 | 32.0873 | −38.8565 | 101.8580 |
| Dividend-to-sales | 4,103 | 1.9161 | 1.1165 | 2.3093 | 0 | 9.2154 |
| IA_Dividend-to-sales | 4,077 | 0.7843 | 0.0701 | 2.2494 | −2.6563 | 7.7169 |
| GS | 4,103 | 6.0426 | 3.4638 | 17.6672 | −77.3508 | 275.0829 |
| GS_decile | 4,103 | 5.5015 | 6 | 2.8725 | 1 | 10 |
| IA_GS | 4,103 | 1.8076 | −0.6801 | 17.6387 | −26.2185 | 270.541 |
| IA_GS_decile | 4,103 | 5.5016 | 6 | 2.8725 | 1 | 10 |
| Tax advantage of dividends | 4,103 | 0.7145 | 0.6670 | 0.1554 | 0.4000 | 1.0750 |

## REFERENCES

Aharony, Joseph, and Itzhak Swary, 1980, Quarterly dividend and earnings announcements and stockholder returns: An empirical analysis, *Journal of Finance* 35, 1–12.

Allen, Franklin, and Roni Michaely, 1997, Dividend policy, in Robert Jarrow, Vojislav Maksimovic, and William Ziemba, eds.: *North-Holland Handbooks in Operations Research and Management Science* (Finance, North-Holland, Amsterdam).

Ambarish, Ramasastry, Kose John, and Joseph Williams, 1987, Efficient signalling with dividends and investments, *Journal of Finance* 42, 321–343.

Asquith, Paul, and David Mullins, 1983, The impact of initiating dividend payments on shareholders' wealth, *Journal of Business* 56, 77–96.

Auerbach, Alan, 1979, Wealth maximization and the cost of capital, *Quarterly Journal of Economics* 93, 433–446.

Baumol, William, 1959, *Business Behavior, Value and Growth* (Macmillan, New York).

Benartzi, Shlomo, Roni Michaely, and Richard Thaler, 1997, Do changes in dividends signal the future or the past?, *Journal of Finance* 52, 1007–1034.

Berle, Adolf, and Gardiner Means, 1932, *The Modern Corporation and Private Property* (Macmillan, New York).

Bhattacharya, Sudipto, 1979, Imperfect information, dividend policy, and the "bird-in-hand" fallacy, *Bell Journal of Economics* 10, 259–270.

Black, Fischer, 1976, The dividend puzzle, *Journal of Portfolio Management* 2, 5–8.

Bulow, Jeremy, and Kenneth Rogoff, 1989, Sovereign debt: Is to forgive to forget?, *American Economic Review* 79, 43–50.

De Angelo, Harry, Linda De Angelo, and Douglas Skinner, 1996, Reversal of fortune: Dividend policy and the disappearance of sustained earnings growth, *Journal of Financial Economics* 40, 341–371.

Demirguc-Kunt, Asli, and Vojislav Maksimovic, 1998, Law, finance, and firm growth, *Journal of Finance* 53, 2107–2138.

Dewenter, Kathryn L., and Vincent A. Warther, 1998, Dividends, asymmetric information, and agency conflicts: Evidence from a comparison of the dividend policies of Japanese and U.S. firms, *Journal of Finance* 53, 879–904.

Easterbrook, Frank, 1984, Two agency cost explanations of dividends, *American Economic Review* 74, 650–659.

Feenberg, Daniel, 1981, Does the investment interest limitation explain the existence of dividends?, *Journal of Financial Economics* 9, 265–269.

Fluck, Zsuzsanna, 1998, Optimal financial contracts: Debt versus outside equity, *Review of Financial Studies* 11, 383–418.

Fluck, Zsuzsanna, 1999, The dynamics of the management-shareholder conflict, *Review of Financial Studies* 12, 347–377.

Gomes, Armando, 2000, Going public with asymmetric information, agency costs, and dynamic trading, *Journal of Finance*, forthcoming.

Harris, Trevor, Glenn Hubbard, and Deen Kemsley, 1997, Are dividend taxes and tax imputation credits capitalized in share values?, manuscript, Columbia University.

Hart, Oliver, and John Moore, 1994, A theory of debt based on inalienability of human capital, *Quarterly Journal of Economics* 109, 841–880.

Jensen, Michael, 1986, Agency cost of free cash flow, corporate finance, and takeovers, *American Economic Review Papers and Proceedings* 76, 323–329.

Jensen, Michael, and William Meckling, 1976, Theory of the firm: Managerial behavior, agency costs, and capital structure, *Journal of Financial Economics* 3, 305–360.

John, Kose, and Joseph Williams, 1985, Dividends, dilution, and taxes: A signalling equilibrium, *Journal of Finance* 40, 1053–1070.

Kang, Jun-Koo, and René M. Stulz, 1996, How different is Japanese corporate finance? An investigation of the information conflict of new security issues, *Review of Financial Studies* 9, 109–139.

King, Mervyn, 1977, *Public Policy and the Corporation* (Chapman and Hall, London).

La Porta, Rafael, Florencio Lopez-de-Silanes, and Andrei Shleifer, 1999, Corporate ownership around the world, *Journal of Finance* 54, 471–517.

La Porta, Rafael, Florencio Lopez-de-Silanes, Andrei Shleifer, and Robert W. Vishny, 1997, Legal determinants of external finance, *Journal of Finance* 52, 1131–1150.

La Porta, Rafael, Florencio Lopez-de-Silanes, Andrei Shleifer, and Robert W. Vishny, 1998, Law and finance, *Journal of Political Economy* 106, 1113–1155.

Levine, Ross, and Sara Zervos, 1998, Stock markets, banks, and economic growth, *American Economic Review* 88, 537–558.

Lintner, John, 1956, Distribution of income of corporations among dividends, retained earnings, and taxes, *American Economic Review* 46, 97–113.

Miller, Merton, and Franco Modigliani, 1961, Dividend policy, growth, and the valuation of shares, *Journal of Business* 34, 411–433.

Miller, Merton, and Kevin Rock, 1985, Dividend policy under asymmetric information, *Journal of Finance* 40, 1031–1051.

Miller, Merton, and Myron Scholes, 1978, Dividends and taxes, *Journal of Financial Economics* 6, 333–364.

Modigliani, Franco, and Merton Miller, 1958, The cost of capital, corporation finance, and the theory of investment, *American Economic Review* 48, 261–297.

Myers, Stewart, 1998, Outside equity financing, Working paper, MIT.

Poterba, James, 1987, Tax policy and corporate savings, *Brookings Papers on Economic Activity* 2, 455–503.

Poterba, James, and Lawrence Summers, 1984, New evidence that taxes affect the valuation of dividends, *Journal of Finance* 39, 1397–1415.

Poterba, James, and Lawrence Summers, 1985, The economic effects of dividend taxation, in Edward Altman and Marti Subramanyam, eds.: *Recent Advances in Corporate Finance* (Richard D. Irwin Publishers, Homewood, Ill.).

Rajan, Raghuram, and Luigi Zingales, 1995, What do we know about capital structure? Some evidence from the international data, *Journal of Finance* 50, 1421–1460.

Shleifer, Andrei, and Robert W. Vishny, 1997, A survey of corporate governance, *Journal of Finance* 52, 737–783.

Zwiebel, Jeffrey, 1996, Dynamic capital structure under managerial entrenchment, *American Economic Review* 86, 1197–1215.

# Agency Theory in the Not-for-Profit Sector: Its Role at Independent Colleges

David E. Olson
*California State University at Bakersfield*

*Agency theory has long been studied in the corporate setting and used to explain performance in management and in boards of directors. However, little has been done to extend this research into the area of not-for-profits. Using data collected from member institutions of the Council of Independent Colleges, relationships between boards of trustees and presidential demography and institutional performance were examined. Data were analyzed using panel regression with a separate panel for each year's data and for each of the responding schools. Using revenue and gift income as dependent variables, it was found that increases in the size, average tenure, and level of business executive background on a board led to subsequent increases in performance for the institution. Diversity of the board had mixed results, whereas presidential tenure improved performance. These findings partially support the hypotheses and extend the explanatory reach of agency theory into the not-for-profit sector.*

## INTRODUCTION

In its broadest sense, an agency relationship exists whenever one person or entity does something on behalf of another. The one taking the action can be referred to as the *agent*, whereas the one it is being done for is known as the *principal*. As society has evolved, this sort of relationship has developed as a common mechanism for using the benefits of expertise and specialization both in business and elsewhere. The focus of agency theory is to study these relationships to determine the most efficient contract between the agents and the principals given that all individuals are self-serving and boundedly rational.

To date, agency theory has been studied in a broad array of business disciplines including economics, finance, and organizational behavior (Eisenhardt, 1989). Most of this research focuses on the relationship between shareholders of publicly traded companies and the managers who work for them.

Although this line of inquiry has proved useful, there is reason to extend the agency research into new venues. In Jensen and Meckling's (1976) seminal work on the theory of the firm, they wrote that "the problem of inducing an agent to behave as if he were maximizing the principal's welfare is quite general. It exists in all organizations and in all cooperative efforts—at every level of management in firms, in universities" (p. 309). This suggests that this theory can and should be extended to view the principal-agent relationship in not-for-profit organizations. Even more fundamentally, Clark (1985) suggested that managers are not agents of shareholders of their respective firms, but rather, are agents of the corporation itself. This legally based perspective takes away the focus of the shareholder, thereby making it easier to apply agency theory in the not-for-profit sector. In both sectors, the board of directors has the power and duty to oversee the organization.

Independent universities and colleges have faced significant challenges during recent years. They have found themselves in an increasingly competitive environment where there are fewer traditional students (Ford, 1990) available to attend all institutions. In a sign of the increase in competition, during the 1960s, roughly half of all undergraduates attended private colleges and universities. In 1992, that number had dropped to 17% (Lord, 1995).

At the same time, government assistance to students through grants has been sharply reduced. Furthermore, other environmental changes such as the need to provide costly technological support for the education of the students and, for some schools, the expansion of a junior college system or a for-profit institution into formerly protected territory have placed a burden on the leadership of independent colleges. For many of these schools, the challenges are significant enough to threaten their existence.

The intent of this study is to test the reach of agency theory research in the area of independent universities and colleges by investigating the relationships between changes in the demographics of their governance structures and institution performance.

THEORY

In its most general form, an agency relationship occurs whenever one individual depends on or engages another to perform some service. In such a relationship, the doer is known as the agent, whereas the affected party is called the principal. Given that the agent is a utility maximizer, is granted decision-making authority (Fama, 1980), and that there are asymmetric levels of information between the two parties (Eisenhardt, 1989), there is reason to believe that the agent will not always act in the best interest of the principal (Berle & Means, 1932).

According to traditional agency theory research, the association between the owners and managers in an open corporation fits this description of an agency relationship. The owners, also known as shareholders, put up the

capital necessary to fund the organization. In return, they receive the residual claims, or profits, that remain after all other claimants are paid. The shareholders are also unrestricted in the sense that they do not need to be involved in the organization in any other way. They are specialists, then, in accepting the risk of capital loss in return for the potential of gain from residual profits. Managers, meanwhile, perform the decision-making function for the organization. They develop the strategic plan for the firm and determine the best method of implementation. In a complex organization, success requires specialization of the decision-making process apart from ownership and also by function, such as management, marketing, and finance. Those with valuable training and knowledge in a particular area are given the responsibility to make decisions in that area, and the decision-making system becomes complicated and hierarchical.

Berle and Means (1932) were among the first to theorize about the characteristics of this relationship. As the capital demands of the firm increase, they require a greater number of owners to fund the expansion. As ownership becomes more diffuse, each owner owns less of the total firm, so any gains or losses have less impact on any given owner. In this way, the motivation to be involved decreases and the owners become increasingly passive. At the same time, specialization increases and the owners lose the ability to understand what sort of decisions should be made and whether management has made good decisions. The power lost by the owners is transferred to the professional managers who determine corporate strategy and implementation. The interests of these managers often diverge from those of the owners (Jensen & Meckling, 1976). Managers maximize a utility function based on compensation, power, security, and status as its central elements (Marris, 1964; Williamson, 1964). Owners are interested in maximizing efficiency and profits (Hill & Snell, 1989). When considering risky but potentially very rewarding projects, managers fully participate in bearing the risk of failure but may receive little or no gain if the project proves successful (Fama & Jensen, 1983). Their perspective, then, is focused more on short-term results and the status quo. Owners, meanwhile, are the beneficiaries of the gains from successful ventures and are then willing to accept a riskier and longer term position for the firm. The agency problems stemming from these divergent utility functions and exacerbated by the dispersion of both owners and management can be reduced by separating the control or monitoring function of the decision-making process from the implementation function. This separation necessitates the addition of a third group so that the same people are not responsible for both making the decisions and evaluating them. In the corporation, this third group is the board of directors. As a body, it is their responsibility to perform the internal control function of the organization.

The actual contractual relationship of a corporation deviates somewhat from this model. By law, the managers of a corporation work not for the

shareholders directly but for the organization itself (Clark, 1985). Whereas the shareholders typically have one objective, to maximize personal wealth, the organization may have many claimants with differing objectives. The board of directors is the ultimate decision-making body of the organization. Therefore, managers are agents first to the organization itself and then to the board. The board, meanwhile, acts on behalf of both the corporation and its shareholders. As they are acting to serve the interests of the shareholders, they too are subject to self-service and bounded rationality (Olson, Koput, Staw, & Barsade, 1996).

NOT-FOR-PROFIT ORGANIZATIONS

In a not-for-profit organization, there are no residual claims to be paid out and no owners expecting to earn a profit. Thus, within these organizations, any conceivable agency relationship between owners and managers is clouded. Furthermore, without residual claims or stock, there is no need for management to worry about the organization being bought or sold in the marketplace. These conditions may suggest that managers in a not-for-profit organization have increased opportunity to pursue self-interest (Dyl, Frant, & Stephenson, 1996). In the place of these owners are the donors of the organization. They contribute to the organization with the expectation that something good will result such as lives being saved, the environment cleaned up, or people educated. Although it is not financial, they anticipate a return from their investment and will invest elsewhere if their expectations are not met. Often, the largest donors also become board members of the organization.

Although the function of a not-for-profit board is similar to for profits, there are some differences that are a result of the absence of residual claims. For example, in for profits, the threat of outside takeover provides the discipline to allow insiders to play a significant role on the board. Without this threat and to prevent collusion or expropriation of funds, not-for-profit boards should be dominated by outsiders (Fama & Jensen, 1983). Furthermore, not-for-profit board members are often substantial donors who serve without pay. Because this shows their interest in the well-being of the organization, it may be assumed that they will take their decision control task seriously. This is particularly important within boards of colleges and universities. As put by the board member of one university, "the first item on every [board meeting] agenda should be whether to fire the president" (Corson, 1960). The decision control role of not-for-profit boards, then, is the same as the for profits. In general, not-for-profit boards also have a special responsibility for generating and managing financial resources. They are often called on to personally contribute to the institution, lead campaigns to encourage others to contribute, and manage the financial resources held by the institution (Rauh, 1969). Together, these two responsibilities, decision control and financial management, are among the most important duties of the boards of private universities and colleges.

## BOARD DEMOGRAPHY

Hambrick and Mason (1984) contended that experience and values, and therefore performance, could be inferred from the demographic characteristics of the members of the top management team in an organization. Since then, demographics has been used to study and predict organizational innovation (Bantel & Jackson, 1989), strategic change (Wiersema & Bantel, 1992), and organizational performance (Hambrick & D'Aveni, 1992; Michel & Hambrick, 1992). Based on this rich history of relationships, I use the demographic profile of the boards to predict performance in this study.

## SIZE

Among private colleges and universities, the size of boards can vary from an average of five members at Catholic institutions to an average of 27 members at other church-related institutions (Rauh, 1969). It is presumed, then, that the size of private college boards varies greatly. Little, however, has been written on the effectiveness of these groups based on their size. From an agency perspective, larger groups have more cognitive resources and knowledge (Bantel & Jackson, 1989; Hambrick & D'Aveni, 1992) and access to more information sources and resources (Hambrick & Mason, 1984), resulting in a larger repertoire of possible practices and greater adaptability (Katz, 1982). As a result, larger boards should have a greater monitoring capacity (Murray, 1989) and be able to access more resources for their organization.

*Hypothesis 1:* Increased board size should increase total revenue and gift income.

## TENURE

Greater board tenure should also increase the board's ability to use information and make it less likely to be persuaded by self-interested arguments of managers. As individuals stay in an organization, they become increasingly confident that they know how to do things the right way (Wanous, 1980). Hence, a board that has been in place longer is likely to be a stronger, more productive board. Furthermore, Olson et al. (1996) found that long-tenured boards were subject to increased levels of escalation of commitment (Staw, 1976). In the not-for-profit sector, board members are generally expected to be personally committed to the success of the organization (Fama & Jensen, 1983), and in private universities, this commitment includes financial contributions to the institution (Ingram & Associates, 1980).

*Hypothesis 2:* Increased board tenure should increase total revenue and gift income.

Arrow (1985) explained what he called "hidden action" as the process of the agent using his or her comparatively superior access to knowledge and power to hide poor performance and a lack of effort from the principals. In for-profit institutions, the existence of a relatively long-tenured chief executive can provide more power in the hands of management versus the board. For instance, Olson et al. (1996) found that long-tenured top management teams in commercial banks were more strategically persistent in their handling of bad loans than were short-tenured teams. More generally, in for profits at least, long-tenured chief executives may have more influence on the selection and retention of the members who sit on the board. Such influence may diminish the board's ability to effectively fulfill their roles of monitoring and controlling management and the organization as a whole and thereby allow management to shirk on their duties and hide their performance. Whereas not-for-profit boards are generally self-perpetuating, the long-tenured college president may still influence the selection of new board members and may be more adept at covering up poor performance. Because a major task of college presidents is to bring in gifts and revenue to the institution, they may use their superior power to shirk on their duties and hide poor performance.

*Hypothesis 3a:* Increased presidential tenure should decrease total revenue and gift income.

However, the roles of management and the board and the relationship between them are different in not for profits than in for-profit organizations. For instance, boards of private universities are generally self-perpetuating and contain few, if any, insiders. This encourages a board that is independent from the president, reducing the potential of board weakness when facing a long-tenured president. With the assurance of a strong board, a long-tenured president may develop an increasing sense of psychological ownership in the success of their institution. As with management ownership in for profits (Olson & Koput, 1998), it is predicted that this increased sense of ownership will lead to increases in performance for the organization.

*Hypothesis 3b:* Increased presidential tenure should increase total revenue and gift income.

BUSINESS EXECUTIVE BACKGROUND

Boards of trustees at private universities are made up of individuals who have backgrounds in many areas such as education, religion, and law. Many trustees report that they are ill prepared for the financial role that they ought to

play for their institution (Ford, 1990). Some, however, have education and experience in organizational decision making, corporate strategy, fund-raising, and financial management. Such experience gives them exposure to making difficult and complex managerial and financial decisions and better understanding and access to financial markets and resources. Such experience has lead to the financial success in for-profit institutions (Boeker & Goodstein, 1991), and at Princeton University, a board member with professional investment experience is credited with making decisions that led to the relatively large size of the school's endowment (Bowen, 1994). Therefore, it is to be expected that such experience will lead to improved performance for their institutions.

*Hypothesis 4:* Board business executive experience should increase total revenue.

## HOMOGENEITY

Homogeneous teams have been seen as more efficient (Hambrick & Mason, 1984; Murray, 1989) because members of homogeneous teams know what to expect from one another (Pfeffer, 1983) and do not have the problems associated with dissimilar experiences, backgrounds, beliefs, and values (Wiersema & Bantel, 1992). Homogeneous background experience was found to lead to better performance on tasks that do not require creativity or innovation (Ancona & Caldwell, 1992), such as financial performance in high-tech firms and lower loan losses in banks (Olson et al., 1996). Furthermore, homogeneous teams tend to communicate more (Murray, 1989) and to be higher in social cohesion (Lott & Lott, 1965). Combining these findings, a more homogeneous board is likely to be stronger and more adept at decision control.

*Hypothesis 5a:* Increased board homogeneity should increase total revenue.

On the other hand, heterogeneous groups have a greater variety of sources from which to gather information (Hambrick & Mason, 1984) and, presumably, fund-raising sources. This variety can lead to greater diversity and comprehensiveness in the set of recommended solutions to a problem. Heterogeneous educational specialization was found to lead to strategic change (Wiersema & Bantel, 1992) and greater political activity (Pfeffer, 1981). When reaching out in fund-raising activities, these attributes may lead to favorable institutional results.

*Hypothesis 5b:* Increased board heterogeneity should increase gift income.

347

METHOD

DATA SOURCES

The data for this research were collected from each participating institution within the Council of Independent Colleges (CIC). This is an administrative body that attracts as members independent colleges and universities from across the United States. The duty of the CIC is to provide support and direction for member institutions and act as a clearinghouse for relevant research studies and political action. To collect the information, a preliminary questionnaire was sent to the CIC for review. After a first round of revisions, the questionnaire was sent to a test school that was not a part of the CIC but whose leadership has been involved in board research at the university level. The feedback from this test was given to the CIC along with a second draft of the questionnaire. The board of the CIC then conducted a multiple school pretest of the questionnaire. From the pretest, final changes were made to the questionnaire before mailing it to the president of each member institution. Accompanying the questionnaire was a letter of endorsement from the president of the CIC and a postage-paid return envelope. Also included was a letter of introduction that explained the nature of the study and assured both anonymity and confidentiality. The survey itself was partitioned into four sections so that each institution could collect and organize responses accurately and efficiently by those with the best access to the information. The four sections were based on institution financial information, student demographics and school information, board and president characteristics, and individual board demographics, respectively. Each section requested information for the five contiguous years from 1991 to 1995. Of the 420 institutions, 43, or 10.2%, returned completed and usable questionnaires (see Tables 1 and 2).

DEPENDENT VARIABLES: FINANCIAL PERFORMANCE

One of the primary duties of the board of an educational institution is managing its financial performance. According to Ingram & Associates, "trustees are in a better position than any other group to preserve and improve the financial health of the institution," and Rauh (1969), "the board of a private college carries the primary responsibility for financing its operation is widely held." Ingram went on to propose that the top guideline for university board effectiveness is to "legitimize the program for obtaining resources." Securing financial resources, then, is a primary objective of boards and a reasonable measure of their effectiveness.

Two financial measures, total revenue and gift income, are used as performance indicators.

**Table 1.   Descriptive Statistics:**
**Characteristics of Schools Boards and Institutional Performance**

|  | Mean | Standard Deviation | Minimum | Maximum | N |
|---|---|---|---|---|---|
| Board size | 28.92 | 7.45 | 11 | 47 | 211 |
| Board tenure | 6.40 | 3.17 | .24 | 15.24 | 215 |
| Percentage with business |  |  |  |  |  |
|   executive experience | .2397 | .1418 | 0 | .62 | 210 |
| Presidential tenure | 8.61 | 6.19 | 0 | 30 | 197 |
| Board functional homogeneity | .2035 | .0707 | .12 | .63 | 210 |
| Board ethnic homogeneity | .8966 | .0969 | .58 | 1.00 | 212 |
| Total revenue (in 1,000s) | 17,876 | 9,415 | 4,027 | 65,652 | 199 |
| Total gift income (in 1,000s) | 2,599 | 1,977 | 192 | 15,117 | 184 |
| Endowment gift (in 1,000s) | 803 | 2,443 | –5 | 29,237 | 184 |
| Total number of gifts | 4,242 | 2,487 | 317 | 13,070 | 165 |

*Note:* Forty-three schools observed over a 5-year period (1991 to 1995). All measures are taken over all institution-years except where data were missing from the institutional response.

**Table 2.   Within-School Correlations Among Characteristics of Boards and Presidents**

| Variable | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 1. Board size |  |  |  |  |  |
| 2. Board tenure | .233 |  |  |  |  |
| 3. Board business experience | .245 | .058 |  |  |  |
| 4. Presidential tenure | –.013 | –.062 | .037 |  |  |
| 5. Homogeneity-function | –.419 | –.265 | .151 | .133 |  |
| 6. Homogeneity-ethnicity | –.095 | .206 | –.001 | –.186 | .170 |

Total revenue is a very broad measure. It is calculated by summing all of the money brought in to the institution during the academic year and includes tuition and student fees, government grants, and other sources of revenue.

Gift income is a narrower measure and is based on the amount of money collected through donations in a given year. Board members are often expected to contribute funds to their institution (Fama & Jensen, 1983). At private colleges and universities, nearly 90% of all responding board members said they were involved in the fund-raising plans of their institutions (Rauh, 1980), and about a third of all board members personally contribute to the capital campaigns of their institutions (Radock, 1977). Board members, then, can make personal contributions to their institutions and are an excellent resource for encouraging others to do the same.

INDEPENDENT VARIABLES: BOARD STRENGTH

The agency theory literature on boards of directors uses a number of demographic measures to estimate board strength. Several are borrowed to test for board performance and strength here. Size is a simple count of the number of directors on the institution's board in a given year. Tenure is computed by taking the numerical average of years served on the school's board for all directors serving in a given year. Business executive experience is computed as the percentage of directors on the institution's board that are business executives in a given year.

### Homogeneity

Because homogeneity is a multidimensional construct (O'Bannon & Gupta, 1992), two dimensions are examined in this study: variation of functional background and variation of ethnicity. In each case, the respondent institutions provided these characteristics for each of their directors. Functional homogeneity was measured by coding directors into the following categories: clergy, law, government, education, medicine, self-employed, agriculture/farming, homemaker, business-finance, business-senior executive, business-other, other, and retired. Next, an index of functional homogeneity was computed for each institution in each year as follows. For institution $i$ in year $t$, I denote the number of directors with a background in occupation type $j$ as $n_{it,j}$ and the total number of directors aggregated over all occupations ($j = 1 \ldots 13$) as $n_{it}$. The proportion of institution $i$'s directors of background $j$, out of the total board size, is denoted $p_{it,j}$ and given by $p_{it,j} = n_{it,j}/n_{it}$. Each $p_{it,j}$ is squared, and then the sum is taken over all $j$, resulting in the index of diversity, $y_{it}$, so that:

$$y_{it} = \sum\nolimits_{J=1}^{J} p_{it,j}^2.$$

This is equivalent to subtracting from unity the index of heterogeneity popularized by Blau (1977).

Ethnicity was measured by coding directors into the following categories: Caucasian, African American, Asian, Hispanic, Native American, and other. Next, an index of ethnic homogeneity was computed for each institution in each year using the same method as described earlier.

### Statistical Methods

Statistical analyses were conducted on five years of cross-sectional records (from 1991 to 1995) for all 43 colleges and universities in the data set, creating a sample size of 215 panels of data to draw from (in most cases, the actual $n$ was smaller than 215 due to lagged variables and missing data; see Table 3). To test the predictions of the hypotheses, I employed a panel regression estimator.

This procedure allowed me to account for the effect of years, individual schools, and other specified control variables. The selection of this technique was primarily based on the theoretical consideration that effects of demography reside within institutions and occur over time. That is, I am interested in explaining processes (e.g., board strength influencing performance) that occur within institutions over time rather than factors (e.g., age of institution) that might determine which institutions had the highest endowments.

As such, it is undesirable for the estimates to be biased by between-institution variation on variables that cannot be observed. Such unobserved heterogeneity could arise due to differences among institutions in omitted variables that are constant over time, such as different initial conditions. Or, unobserved heterogeneity might result from differences over years in omitted variables that are constant over institutions, such as changes in economic conditions. These omitted variables could affect both independent and dependent variables (as a common cause), biasing estimates of the parameters (capturing the relationship between independent and dependent variables). For example, some schools may have weaker boards and lower revenues due to differences in histories or strategies. To eliminate any spurious effects due to unobserved differences among institutions, a random school effect was included in the model and tested for correlation between the error and the explanatory variables (Hsiao, 1986). That is, I included a parameter to capture the variance between schools. This random effects approach is used rather than the alternative fixed effects specification sometimes used in panel regression because the sample does not include the entire population but does appear representative on most school characteristics such as size and geographic location. Autocorrelation may also bias parameter estimates because of factors that change over time within institutions but are not included in the model. For example, institutions could have cycles of revenue or gifts that have naturally evolving patterns that change in coherent but unforeseeable ways over time. To control for this, I included a lagged dependent variable as a predictor in each model.

## RESULTS

Table 3 contains the results of the panel regression. There are four dependent variables, each capturing an aspect of performance. The model in the first column predicts total revenue for the institutions. This is the sum of all the financial resources collected by the institution during a given year. Each of the remaining three columns represents more specific measures of institutional performance. Columns 2 and 3 represent the amount and number, respectively, of total gifts given to the institution in a given year. These gifts may include annual unrestricted gifts, capital expansion gifts, deferred gifts, and other miscellaneous gifts. Endowment gifts is the dependent variable in the

**Table 3.   Results of Panel Regressions: Effects of Board Demography on School Performance**

| Independent Variable | Total Revenue (in $1,000s) | | Total Gifts (in $1,000s) | | Total Gifts (in units) | | Endowment Gifts (in $1,000s) | |
|---|---|---|---|---|---|---|---|---|
| | *Dependent Variable* | | | | | | | |
| Full model | | | | | | | | |
| 1. Board size | .1495 | (.1016) | .0911 | (.0294)*** | .0575 | (.0197)*** | .0500 | (.0155)*** |
| 2. Board tenure | 1.1585 | (.3102)*** | .1554 | (.0771)** | .1646 | (.0718)** | −.0015 | (.0384) |
| 3. Business executive background | 13.8307 | (6.5952)** | .0015 | (1.6951) | 2.031 | (1.3640) | −.0997 | (.8476) |
| 4. Board homogeneity-background | −1.9749 | (9.2482) | 7.9727 | (2.8754)*** | −2.5863 | (1.7431) | 4.5146 | (1.4862)*** |
| 5. Board homogeneity-ethnicity | 3.5245 | (4.1747) | −3.4178 | (1.2829)*** | 1.4890 | (.8648)* | −2.1377 | (.6501)*** |
| 6. Presidential tenure | −.0431 | (.0716) | .0511 | (.0279)* | −.0166 | (.0121) | .0239 | (.0139)* |
| 7. Random effects | .1799 | | .4326 | | .0956 | | .5176 | |
| 8. $R^2$ | .9010 | | .6623 | | .9692 | | .5806 | |
| 9. N | 147 | | 136 | | 124 | | 140 | |

*Note:* Standard errors in parentheses.
*$p < .10$. **$p < .05$. ***$p < .01$.

fourth column. The rows of Table 3 represent each of the predictor variables in the study.

Looking first at row 1, there was found to be a positive relationship between board size and the dependent variables. More specifically, as the number of members who composed the board increased (decreased) in one period, gift income, the total number of gifts, and endowment gifts increased (decreased) in the subsequent year. This provides support for Hypothesis 1.

Row 2 measures the effect of board tenure on the dependent variables. As the average length of service of the board members lengthened (shortened), total revenue, gift revenue, and the total number of gifts given to the institution tended to increase (decrease). This provides support for Hypothesis 2.

The third row looks at the effect of having board members with business executive backgrounds. As the percentage of members that were business executives increased (decreased) in one year, total revenue increased (decreased) in the following period. Business executive background had no significant impact on gift income, the number of gifts received by the school, or endowment gifts. Thus, there is support for Hypothesis 4.

Rows 4 and 5 illustrate the effect of board homogeneity on the dependent variables. Looking at row 4, as a board becomes more homogeneous (heterogeneous) with respect to functional background, total gift income and endowment gifts will increase (decrease) in the subsequent year. However, row 5 illustrates the opposite pattern as the board becomes more homogeneous with respect to ethnicity. In this case, a more ethnically homogeneous (heterogeneous) board will lead to less (more) gift income, total number of gifts, and endowment gifts in the subsequent year of the institution. These results provide mixed and conflicting support for Hypotheses 5a and 5b.

Row 6 measures the effect of presidential tenure on the revenue and gifts. As presidential tenure increased (decreased), the total number of gifts and endowment gifts tended to increase in the subsequent year. These results refute Hypothesis 3a and provide some support for Hypothesis 3b.

## DISCUSSION

As expected, the demographic characteristics of individual board members and of the board as a whole played a role in determining the performance outcomes of the colleges and universities with which they served. Because demography can serve as a proxy for the values and experiences that a person holds (Hambrick & Mason, 1984), we can expect that by looking at selected demographic characteristics of board members, we can gain insight into how that board may perform.

Looking more specifically, it was found that school revenue and gifts were enhanced when board members had relatively longer average tenures. There are several possible explanations for these relationships. One is the standard

agency theory position. Here it is argued that long-tenured boards are less likely to be selected by or feel under the control of the existing management team. This makes them more independent and stronger and therefore able to do a better job of monitoring management. Management, then, being more closely watched and assisted, is able and willing to do a better job of bringing resources into the institution. A variation to this agency explanation is that the relative strength of a long-tenured board causes not just increased monitoring but also a greater sense of responsibility on the part of the board (Olson et al., 1996). Under this scenario, the board itself does more of the work as its tenure increases. An additional variation is that longer tenure may also lead to a greater level of psychological commitment to the institution (Olson et al., 1996). Here, the board members themselves are not just watching and working more, they are also increasing their personal gifts to the schools as they increase their tenure on the board. Finally, it may be theorized that there is a self-selection bias whereby members of successful schools choose to stay at their institutions for a longer period of time. Although this may be true, it does not seem to adequately explain that increases in revenue and gifts are preceded by the increases in tenure.

It was also found that there is a positive relationship between board size and gifts to the institution. As with the original agency argument presented with board tenure earlier, here a larger board may feel stronger and more independent, thereby providing more control over and direction to management. However, because revenue did not rise significantly and gifts did, there is reason to consider a second theory, resource dependence, as a better, or at least additional, explanation of these results. This theory emphasizes that organizations are dependent on external resources to survive and thrive (Pfeffer & Salancik, 1978). When possible, the organization will try to secure these resources through any means possible including cooptation. In this case, the universities may bring in additional board members because they need access to gifts. The theory suggests that the members are selected by their ability to bring in gifts. It is clear from the data that it is not just the new members who are giving gifts (otherwise, there would not be a significant increase in the absolute number of gifts received). However, the new members may well have access to new sources of gift givers, thus explaining both the increase in gift income and in the number of gifts given.

Although business executive background did not lead to more gifts, it is associated with increases in total revenue. This finding is consistent with the agency arguments of a stronger board leading to greater control and access to information.

The results of board homogeneity on school performance were mixed and rather puzzling. First, neither measure of homogeneity had any significant impact on total revenue brought into the institutions. With regard to ethnicity, a more heterogeneous board led to more gift income. As with increases in the size of the board, this result may be better explained using resource

dependence theory. In the majority of cases, the boards were not very ethnically diverse to begin with. By increasing ethnic diversity, then, the institution could avail itself to a new source of givers that may feel more included in the mission of the institution. With regard to functional background, a more homogeneous board led to more gift income. Because homogeneous teams tend to communicate better, this may make them better at the control function. This, then, provides some support for the agency argument. However, because total revenue did not increase along with functional homogeneity, there is reason to suspect there is another unknown cause of this relationship.

Finally, it was found that presidential tenure had a positive impact on gift income at the institution. On the surface, this argues against agency theory, as a long-tenured president gains additional power and control over the institution. However, if we consider liberalizing ownership to include a psychological commitment, a long-tenured president may feel a greater sense of ownership to his or her institution, thereby making them feel more like principals than agents. Still, I concede that this is only one possible explanation for these results.

## CONCLUSION

There are several implications for these results of this research. First, it is important to note that not-for-profit boards can and do make a difference in the institutions they serve. They need not be there just for show or to fulfill regulatory requirements but can make a significant contribution to the performance of the organization. Selection of members, then, becomes an important function. Boards tended to improve performance as they grew larger and longer tenured, and gained members with business executive background. Board diversity garnered mixed results, with ethnic diversity improving performance and functional background diversity hurting performance. Presidential tenure also seemed to improve performance, at least when measured by gift income. These results provide some support for the extension of agency theory to the not-for-profit sector and also for the liberalization of the theory to include psychological ownership. Support for resource dependence theory was also found. Overall, the results suggest that more research can and should be conducted to test agency theory and board performance in the not-for-profit sector and that demography can be used as an effective tool in measuring and enhancing board performance.

References

Ancona, D. G., & Caldwell, D. F. (1992). Demography and design: Predictors of new product team performance. *Organization Science*, *3*, 321-341.

Arrow, K. J. (1985). The economics of agency. In J. W. Pratt & R. J. Zeckhauser (Eds.), *Principals and agents: The structure of business* (pp. 37-51). Boston: Harvard Business School Press.

Bantel, K., & Jackson, S. (1989). Top management and innovations in banking: Does the composition of the top team make a difference? *Strategic Management Journal*, *10*, 107-124.

Berle, A., & Means, G. (1932). *The modern corporation and private property*. New York: Macmillan.

Blau, P. (1977). *Inequality and heterogeneity: A primitive theory of social structure*. New York: Free Press.

Boeker, W., & Goodstein, J. (1991). Organizational performance and adaptation: Effects of environment and performance on changes in board composition. *Academy of Management Journal*, *34*, 805-826.

Bowen, W. G. (1994). *Inside the boardroom*. New York: John Wiley.

Clark, R. C. (1985). The economics of agency. In J. W. Pratt & R. J. Zeckhauser (Eds.), *Principals and agents: The structure of business* (pp. 55-79). Boston: Harvard Business School Press.

Corson, J. J. (1960). *The governance of colleges and universities*. New York: McGraw-Hill.

Dyl, E. A., Frant, H. L., & Stephenson, C. A. (1996, January). *Governance structure and performance of not-for-profit corporations: Evidence from medical research charities*. Paper presented at forum, University of Arizona, Tucson.

Eisenhardt, K. M (1989). Agency theory: An assessment and review. *Academy of Management Review*, *14*, 57-74.

Fama, E. F. (1980). Agency problems and the theory of the firm. *Journal of Political Economy*, *88*, 288-307.

Fama, E. F., & Jenson, M. C. (1983). The separation of ownership and control. *Journal of Law and Economics*, *26*, 301-325.

Ford, G. F. (1990). Trustees: Taking fund raising responsibility. In W. K. Willmer (Ed.), *Friends, funds and freshmen: A manager's guide to Christian college advancement* (pp. 41-54). Washington, DC: Christian College Coalition.

Hambrick, D. C., & D'Aveni, R. A. (1992). Top team deterioration as part of the downward spiral of large corporate bankruptcies. *Management Science*, *38*, 1445-1466.

Hambrick, D. C., & Mason, P. (1984). Upper echelons: The organization as a reflection of its top managers. *Academy of Management Review*, *9*, 193-206.

Hill, C., & Snell, S. (1989). Effects of ownership structure and control on corporate productivity. *Academy of Management Journal*, *32*, 25-46.

Hsiao, C. (1986). *Analysis of panel data*. New York: Cambridge University Press.

Jensen, M. C., & Meckling, W. H. (1976). Theory of the firm: Managerial behavior, agency costs, and ownership structure. *Journal of Financial Economics*, *3*, 305-360.

Katz, R. (1982). The effects of group longevity on project communication and performance. *Administrative Science Quarterly*, *27*, 81-104.

Lord, M. (1995, September 18). A battle for survival. *U.S. News & World Report*, 132-134.

Lott, B. E., & Lott, A. J. (1965). Group cohesiveness and interpersonal attraction: A review of relationships with antecedent and consequent variables. *Psychological Bulletin*, *4*, 259-309.

Marris, R. (1964). *The economic theory of managerial capitalism*. New York: Macmillan.

Michel, J. G. & Hambrick, D. C. (1992). Diversification posture and top management team characteristics. *Academy of Management Journal*, *35*, 9-37.

Murray, A. (1989). Top management group heterogeneity and firm performance. *Strategic Management Journal*, *10*, 125-141.

O'Bannon, D. P., & Gupta, A. K. (1992, August). *Utility of heterogeneity versus homogeneity within top management teams: Towards a resolution of the empirical paradox*. Paper presented at the meeting of the Academy of Management, Las Vegas.

Olson, D., & Koput, K. (1998, August). *For better or for worse: Effects of management stock ownership on performance, commitment, and agency in a longitudinal study of California banks*. Paper presented at the meeting of the Academy of Management, San Diego.

Olson, D., Koput, K., Staw, B., & Barsade, S. (1988, August). *Agency, escalation, and hypocrisy in a longitudinal study of banks' coping with problem loans*. Paper presented at the meeting of the Academy of Management, Cincinnati, OH

Pfeffer, J. (1981). *Power in organizations*. New York: McGraw-Hill.

Pfeffer, J. (1983) Organizational demography. In L. L. Cummings & B. M. Staw (Eds.), *Research in organizational behavior* (Vol. 5, pp. 299-357). Greenwich, CT: JAI.

Pfeffer, J., & Salancik, G. R. (1978). *The external control of organizations*. New York: Harper and Row.

Radock, M. (1977). *The fund-raising role*. Washington, DC: Association of Governing Boards of Universities and Colleges.

Rauh, M. A. (1969). *The trusteeship of colleges and universities*. New York: McGraw-Hill.

Smith, K. G., Smith, K. A., Olian, J. D., Sims, H. P., Jr., O'Bannon, D. P., & Scully, J. A. (1994). Top management team demography and process: The role of social integration and communication. *Administrative Science Quarterly*, *39*, 412-438.

Staw, B. M. (1976). Knee-deep in the big muddy: A study of escalating commitment to a chosen course of action. *Organizational Behavior and Human Performance*, *16*, 27-44.

Wanous, J. (1980). *Organizational entry: Recruitment, selection, and socialization of newcomers*. Reading, MA: Addison-Wesley.

Wiersama, M. F., & Bantel, K. A. (1992). Top management team demography and corporate strategic change. *Academy of Management Journal*, *35*, 91-121.

Williamson, O. E. (1964). *The economics of discretionary behavior: Managerial objectives in a theory of the firm*. Englewood Cliffs, NJ: Prentice Hall.

*David E. Olson is an assistant professor of management at California State University at Bakersfield where he teaches Strategy, Entrepreneurship, and Organizational Behavior. His research interests include governance in both the for-profit and not-for-profit sectors and social dilemmas. He obtained his Ph.D. from the University of Arizona in 1998, an M.S. from the University of Arizona in 1997, an M.B.A. from the University of Washington in 1987, and a B.A. from Westmont College in 1981. He also has extensive experience working with financial institutions in the areas of strategy, financial analysis, and general management.*

# Introduction

# Principal–agent theory and research policy: an introduction

## Dietmar Braun and David H Guston

*The rational choice perspective is prominent in many sociological, economic and political science literature but has been undervalued until now in the field of science studies. This special issue attempts to revalorise this perspective by introducing the principal–agent theory with relation to research policy-making. The introduction presents the basic features of the model of principal–agent and reviews the theoretical development and applications in research policy. It summarises the main findings of the articles in this issue and concludes that the studies in the framework of principal–agent demonstrate the willingness of combining theoretical rigour and 'requisite variety' by applying the theory to a large number of different fields linked to research policy-making.*

Professor Dietmar Braun is at the Institut d'Etudes Politiques et Internationales, Université de Lausanne, BFSH 2, CH-1015 Lausanne, Switzerland; Tel: +41 21 6923132; Fax: +41 21 6923145; E-mail: Dietmar.Braun@iepi.unil.ch. Professor David H Guston is in the Department of Public Policy, Rutgers, The State University of New Jersey, 33 Livingston Ave, Suite 202, New Brunswick, NJ 08901-1980 USA; Tel: +1 732 932 2499 707; Fax: +1 732 932 1107; E-mail: guston@rci.rutgers.edu.

ALTHOUGH THEY HAVE never played a prominent role in the sociology of science or political science literature, studies on research policy-making have a long tradition. Furthermore, it must be stated that a common theoretical framework has been lacking. The field is, as David Guston (1996) noticed, heavily "undertheorised". It would be wrong, however, to assume that there were no attempts in the past to learn lessons about research policies on the base of various theoretical approaches. Though a systematic and historical overview is lacking, we can discern perhaps five theoretical currents that have influenced in some way the thinking on research policies:

- Economics, both classical (Polanyi, 1951; 1962; Tullock, 1966; Ghiselin, 1987; Foray, 2000) and Marxist (Bourdieu, 1975; Bourdieu, 2001);
- System-theory (Krohn and Küppers, 1987; Luhmann, 1990);
- Constructivism (Latour and Woolgar, 1979; Knorr *et al*, 1981; Knorr-Cetina, 1981; Latour, 1987);
- Institutionalism (Merton, 1970 [1938]; Ben-David, 1971; 1991; 1991 [1977]; Mukerji, 1989; Mayntz, 1991);
- Discussion on the "finalization of science" (Daehle, 1979; Weingart, 1997).

This introduction is not the place to assess the contribution of each approach to our present knowledge on research policies. We generally believe, however, that these approaches are either too abstract (system-theory), lacking in parsimony or theoretical rigour (constructivism, institutionalism, finalization) or, if they are parsimonious (like classical economics), that they abstract too greatly from the 'requisite variety' of real life.

Dietmar Braun, professor of comparative political science at the Institut d'Etudes Politiques et Internationales of the University of Lausanne (Switzerland), has studied political sciences at the University of Amsterdam (PhD) before he worked as a research fellow at the Max Planck Institute for Societal Research in Germany. In 1996, he presented his habilitation thesis at the Institute for Political Science at the University of Heidelberg before he became, in October 1996, full professor in Lausanne. Dietmar Braun has published on labour market, higher education and research policies, federalism as well as on modern political theory. He has published several books, and articles in journals.

David H Guston is associate professor and director of the Public Policy Program at the Bloustein School of Planning and Public Policy at Rutgers University. His book, *Between Politics and Science* (Cambridge University Press, 2000) was awarded the 2002 Don K Price Prize by the American Political Science Association for best book in science and technology policy. His published work focuses on research and development policy, scientific integrity and responsibility, science advice and public participation in technical decision making, peer review, and the politics of science policy. His current research includes investigating the public value of social policy research sponsored by US federal social policy agencies.

What we need instead, is an intelligent combination of analytical and rigorous tools, of "parsimony, refinement, and (in the sense used by mathematicians) elegance" (Bates *et al*, 1998, page 11). Such tools will be useful for interpreting research policies with attention to the historical and institutional contexts in which research policy is made, something we find for example in the concept of "analytic narratives" developed by Bates *et al* (1998).

The analytic tool most useful for this purpose seems to us to be principal–agent theory, which has been developed in the context of rational choice and transaction cost theory (see, for example, Ross, 1973; Williamson, 1975; 1985; Coleman, 1990; Kiewiet and McCubbins, 1991). In the early 1990s, Braun (1993) introduced the concept in the context of research policy-making (see below), most notably by referring to the relationship between policy-makers as the principal and the various funding agencies responsible for the implementation of research policy as the agents. This relationship seemed to correspond perfectly to the basic logic of principal–agent figurations, that is, one actor who seeks "extension of self" (Coleman, 1990, page 146) by delegating some tasks for execution by other actors who seem better capable to do so.

Funding agencies were, since their origins, designed to work out and implement research policies, in preference to the usual public bureaucracy that lacked the necessary direct contacts with science. The concept is, however, more general and can, as will be shown in this special issue, be applied to relations between policy-makers and scientists in general, between the funding administration and scientists, or between funding programme directors and scientists.

Principal–agent theory is becoming a predominant approach in different fields of political science where 'delegation' as one particular form of organising state activities is discussed. This predominance holds for studies of delegation to bureaucracy by Congress (McCubbins, 1984; Weingast, 1984), delegation to "independent regulatory agencies" (Majone, 2001a) and central banks (Majone, 2001b; Elgie, 2002), as well as of all relations in the "democratic chain of delegation" from the voter to the implementing bureaucracy (Strom, 2000). Given the extant literature on principal–agent theory in research policy, we consider it appropriate to demonstrate the usefulness of the approach in this policy field and to further elaborate the concept.

We want, first, to briefly introduce the reader to the basics of principal–agent ideas before we then present the findings on research policy in the literature. Finally, we will give an overview of the main themes discussed in this special issue.

## Basics of principal–agent theory

The principal–agent literature deals with a specific social relationship, that is, delegation, in which two actors are involved in an exchange of resources. The principal is the actor who disposes of a number of resources but "not those of the appropriate kind to realize the interests (for example, has money but not the appropriate skills)" (Coleman, 1990, page 146). He or she then needs the agent, who accepts these appropriate resources and is willing to further the interests of the principal. In this sense, Coleman is right to speak of an "extension of self" of the principal by way of delegation.

The principal–agent model has been developed within the framework of the "new institutional economics" (Williamson, 1975; 1985; Moe, 1984; Miller, 1992) and therefore shares the basic characteristics of this framework, for instance, the assumption of rational actors striving to maximise their preferences that are ordered according to their priorities. Institutions can constrain actors' choices so that the conscious design of institutions (like contracts) may help to overcome typical collective action problems involved in the principal–agent relationship.

There are two typical collective action problems discussed in the literature — moral hazard and adverse selection. These problems are based on what the new institutional economics calls the 'opportunism' of actors: Actors are self-interested and thus seek to maximise their personal welfare. They may do this by seeking their self-interest "with guile. This includes but is scarcely limited to more blatant forms, such as lying, stealing, and cheating" (Williamson, 1985, page 47).

In the particular case of principal–agent relations, such 'cheating' or, as is often said, 'shirking' by the agent may happen because the agent usually has an informational advantage *vis-à-vis* the principal. The principal does not know for sure if the agent will really do his or her best when delegated certain tasks

(this is the "moral hazard"), and usually the principal does not have sufficient information on the abilities of potential agents to find the one best suited to do the task (this is "adverse selection").

As agents seek their self-interest with guile, they may hide this information from the principal to reduce their work load or to be hired in the first place. The resulting delegation could then be sub-optimal or even detrimental to what the principal attempted to achieve. This is why the principal–agent literature (discussing the problems of insurance agencies or of parliaments dealing with bureaucracies) discusses contract and monitoring mechanisms designed to avoid these problems.

The collective action problems arise as both sides — the principal and the agent — have an interest in entering into the exchange relationship. They both profit by exchanging resources: the principal by getting something done he or she could not otherwise do, and the agent because he or she gets remuneration of some kind (money, social recognition, and so on). Despite these mutual advantages, the collective outcome may be suboptimal because, as is said, the agent has incentives to seek his self-interest with guile.

We should also not forget the possibility of the principal to 'shirk', a possibility often not discussed in the literature. He or she may have incentives not to deliver the resources fully as agreed to in the contract. Because of these co-operative and selfish motives characterising the relationship, principal–agent interaction is a 'mixed-motive game'.

The principal–agent literature discusses not only contract and monitoring mechanisms but also different possible configurations and their influence on the shirking of agents. A configuration reducing the possibilities to shirk, for example, is the presence of multiple agents (Kiewiet and McCubbins, 1991; Ferejohn, 1993), which creates more of a market-like structure. The more advantageous configuration for the agents of multiple principals has been less often discussed (but see Elgie, 2002).

## Principal–agent theory enters science policy

Principal–agent theory began to enter science policy in the 1990s. Guston (1996) wanted to use it to "re-interpret" generic science policy problems. For him, "the problem of science policy is the problem of delegation" because lack of information on the side of "non-scientists" leads to the typical problems of adverse selection and moral hazard. This formulation provides the opportunity to reflect about adequate incentive structures to solve the main problems in science policy such as the integrity and productivity of research, or the choice between mission and disciplinary research.

The main thrust of his paper was to make understood that treating scientists as agents does not at all mean a hierarchical relationship. The autonomy of

> **A configuration reducing the possibilities of agents shirking is the presence of multiple agents, which creates more of a market-like structure: the more advantageous configuration for the agents of multiple principals has been less often discussed**

agents is widely respected in this relationship, and one should consider the relation between policy-makers and scientists as a "two-way street" where a certain degree of autonomy is respected on both sides.

Neither does Guston want to defend the usual normative stance of the principal–agent literature, which is to inform the principal about how best to design incentive structures. Principal–agent can be used for a variety of purposes and illuminate the basic tensions in the generic problems mentioned above, and it can be used to reflect in general on the "fairness" of the contracts defined between the principal and the agent.

Guston applied principal–agent theory to the relationship between non-scientists and scientists in general. Most of the other literature — and Guston (2000) joins them later — deals with a very specific feature of science policy, that is, the use of research councils and funding agencies that intermediate between policy-makers on the one hand and scientific agencies and scientists on the other hand. The introduction of an intermediary level makes the discussion more complex. The central question becomes: in what way can research councils be seen as agencies that serve the interests of policy-makers?

The introduction of research councils can, from a functionalist perspective, be explained by the proximity of these agencies with the scientific field, more information on the field, and a greater capacity to aggregate the available knowledge. From an economic perspective, the establishment of funding agencies can be explained by the decrease in transaction costs for policy-makers in developing science policies. It seems, in addition, easier to influence such agencies that are either public or semi-public than to influence the scientists themselves because of the former's direct dependence and constitutional commitments to work in the interests of policy-makers.

Principal–agent theory draws the attention of the observer immediately to the possibility of 'shirking' by research councils as agents. Since the beginning, however, contributions in science policy have explained such shirking not simply in terms of

'opportunistic' behaviour but in relation to the interaction between funding agencies and scientists. Braun (1993) criticised the usual dyadic way of conceptualising principal–agent relationships and offered a theoretical account based on the "triadic relationship" among policy-makers, funding agencies and scientists.

His main argument is that funding agencies become intimately interwoven with the "third party" — an under-developed concept in principal–agent theory but already mentioned in the work of Coleman (1990) — to fulfil the task delegated by the policy-maker principal. We could conceptualise the triadic relationship by designing two separate principal–agent relationships, one between policy-makers and the research council and one between the research council and scientists (see also Rip, 1994). Braun would find this too simple, however.

First, scientists are certainly not simply agents of research councils but have a high degree of autonomy and influence on what is decided in research councils. The same holds for the relationship between policy-makers and research councils. We should rather speak of interdependent relationships in which both sides have something necessary to offer the other and a certain degree of autonomy is crucial for all actors in the game. This is, in fact, the "two way relationship" Rip and van der Meulen (1996) evoke, only that we now have *two* 'two way relationships'.

Second, each relationship is influenced by the way the relationship on the other side is organised. If funding agencies choose to 'shirk' in favour of the interests of scientists, this behaviour will have repercussions for the strategies of principals to organise the principal–agent relationship with research councils. If policy-makers are using their formal authority to oblige funding agencies to comply with political interests, this will have implications for the way scientists co-operate with research councils. In such dual interdependent relationships, research councils fare best when they are able to balance the often opposed interests of scientists and policy-makers. Such balancing demands a considerable degree of independence with respect to the principal and to the third party (see also Rip, 1994, page 13; Caswill, 1998).

Both Rip and Braun underline, in addition, the historical and institutional context that is decisive for how the triadic relationship has evolved. There were times, especially during the "science-push period" after World War II, when scientific interests became predominant and research councils seemed to be 'captured' by these interests. Today we seem to experience a period when research councils are more and more captured by political interests. Both periods have, or will, in the end, destabilise the triangle and will lead to increasing efforts of funding agencies to find a new equilibrium.

This idea of finding 'stable equilibria' in science policy principal–agent relationships was put forward in an analytical innovative way, by using game theory (which has found some attention recently in principal–agent literature (Huber and Lupia, 2001)) by van der Meulen (1998). He does not deal with the triadic relationship but instead discusses the "basic" relationship between policy-makers and scientists.

He demonstrates two things: first, that given the utility functions of both the principal and the agent and their options (the principal: to trust or to monitor; the agent: to comply or not comply), the game does not find a stable Nash equilibrium. Each choice creates an incentive for at least one of the two actors to change the *status quo*. Nevertheless, and this is the second point, there are possibilities to find stable equilibria if we assume that, because of the interdependency of actors, the relationship is a long-term one and the game can become co-operative.

Van der Meulen sees four stabilising structures: the role of funding agencies, as discussed before, that can serve a balancing function; the mutual interest of scientists and policy-makers in using peer review; the emergence of a consensus on policy goals; and a competition among agents organised by the principal. For each structure, van der Meulen finds that scientists and policy-makers can develop an interest in co-operation instead of opportunism.

Finally, Guston (2000) also discusses stabilising arrangements in science policy, but finds them above all in the existence of "boundary organisations". His reasoning is influenced by the constructivist literature that he sees as a useful corrective to the often "stylistic" assumption of principal–agent theory. Within the constructivist approach, "boundary" is a term developed by Gieryn (1995) to characterise the often fluid and ambiguous demarcations between scientific and non-scientific fields.

Guston builds his boundary organisations on the notion of "boundary objects" that "allow members of different communities to work together around them, and yet maintain their disparate identities" (Guston, 2000, page 29). Boundary organisations can then be seen as situated between politics and science, both of which can be regarded as principals to the boundary organisations, and, "in doing so, [the boundary organisations] internalise the provisional and ambiguous character of that boundary". They fulfil, therefore, exactly the stabilising function that van der Meulen ascribes to funding agencies, and which is also alluded to in the texts of Braun and Rip.

Boundary organisations may also be of another kind: Guston (2000) treats the examples of the Office of Research Integrity and the Office of Technology Transfer in the United States. The most important point of his study is perhaps that he underlines — and this is the stabilising function of boundary organisations — the inherent capacity of these organisations to facilitate "co-production", that is, the creation of both "knowledge and social order" or, in other terms, both scientific and political interests. Because boundary organisations internalise the

different logics of action, they can bridge different "worlds". The result of his study confirms that intermediary organisations and boundary organisations are crucial to stabilise, according to the logic of principal–agent theory, the inherent unstable relationship of politics and science.

The application of the principal–agent approach in this 'first wave' of science policy studies demonstrates that the theory is not applied without reflection and modifications. The 'field', in particular the position of intermediary organisations, provoked a more complex reflection, which led to the introduction of 'triadic relationships', 'boundary organisations', 'equilibria in games', and 'dynamics in interdependent relationships'.

## Principal–agent in science policy today

This special issue unites most of the authors who have participated in the 'first wave' of publications and adds some others. The articles presented here demonstrate the attempt to broaden the field of application of principal–agent theory in science policy studies and brighten the interesting light this perspective can throw on the choices and procedures of science policy actors.

Braun elaborates the existing insight that historical periods in the funding of science create different utility functions for both political principals and scientific agents, which lead to various "games". The different periods of science policy-making after World War II are interpreted from the angle of principal–agent to see in what way the basic antinomy of funding policies (that is, the maintenance of the autonomy of scientists and the political interest to influence scientific action) is treated within these periods. His article ends by pointing to two different ways in today's research policies of organising the principal–agent relationship, either in a market-oriented way, which increases the moral hazard of scientists, or by funding inter-systemic networks, which seems to be a promising way to overcome the problem of moral hazard and monitoring costs in science policy.

Van der Meulen takes up the classical issue of the role of funding agencies. He increases, however, the complexity of the configurations these agencies have to deal with because recent developments in science policy have shown that users are becoming a kind of "fourth party" in the principal–agent game. By using an empirical study of the Norwegian Research Council, which comprises all types of funding under one roof (but in different divisions), he looks for the strategies that different divisions develop given the various configurations of actors they are dealing with. He finds that differences with respect to the acceptance of strategic funding given the contacts and relationships with users are indeed quite important. Actor constellations, understood as interdependent and multiple relations, matter for what

> **The articles presented here demonstrate the attempt to broaden the field of application of principal–agent theory in science policy studies and brighten the interesting light this perspective can throw on the choices and procedures of science policy actors**

funding agencies are doing. Principal–agent relationships must therefore be understood in the light of these multiple and interdependent configurations, thus reinforcing Braun's finding.

Caswill takes seriously the idea of the contract that figures so prominently in principal–agent literature and embarks on a discussion of how the principal–agent relationship between funding agencies and scientists is organised by actual funding contracts. On the base of an exploratory, empirical study, he shows variations and similarities in how funding agencies set up contracts and attempt to monitor them. He confirms that funding agencies play a crucial role in "mediating" the harsh exigencies in the principal–agent relationship between politics and science. It turns out that the contracts are rarely monitored and that scientists have a considerable freedom to deal with these contracts.

On the other hand, one finds almost no 'shirking' by scientists, meaning there is seldom abuse of this freedom. This also confirms the interdependent relationship between funding agencies and scientists: above all, Caswill highlights, the scientific staff in funding agencies has an intrinsic interest in granting scientists sufficient freedom for action. Scientists, on the other hand, have an interest in complying with the funding agency to be sure of future resources.

Both Guston and Morris attempt to turn the usual top-down perspective of the principal–agent discussion around and take in a bottom-up perspective.

Guston analyses scientific advice as a form of "science in policy". In fact, one of the main problems for policy-makers is to know, given the often contradictory or competing advice different scientists may give, which advice to trust. This is, indeed, the problem of "adverse selection", which has rarely been treated in the context of science policy. He demonstrates that mediation and procedures as well as the creation of market-regulating mechanisms can be an effective way to overcome this problem.

Guston's study is, therefore, very much a debate on the "institutional design" of effective policy-making, only this time seen in the context of policy-formulation. He points to the fact that principal–agent is not a complete theory for scientific advice as such, but that it is very helpful in understanding

the choices and procedures of organising the use of reliable knowledge in policy contexts.

Morris looks at the scientists, the agents in principal–agent theory, and wonders in what way they really behave as agents and how they experience the relationship with policy-makers as principals. By highlighting several "contextual features" of the scientific system, she is able to demonstrate that these features mitigate the sharp edges of the principal–agent relationship. Because of these features — among them, again, is the role of funding agencies — scientists have more freedom than we might expect, there are fewer conflicts between scientists and policy-makers, and accountability is less demanding than we might believe. These structural features allow scientists to do their work without feeling that they are merely agents.

Shove, finally, is the most critical with respect to principal–agent theory, because the experience with several research programmes has demonstrated that the means of principals' influence over scientists through research programmes is extremely limited. Rather, one finds an "anarchic structuring of the field". Shove conceptualises research programmes as the relationship among the funding agency, programme directors, and scientists as multiple agents. Neither funding agencies nor programme directors are able to inhibit the development of a dynamic of research programmes where scientists use these programmes for their own purposes, build up networks, and participate in a multitude of different programmes. These programmes are for her not agents, as she has presumed in the beginning, but they become their own actors. Principal–agent theory cannot, according to Shove, understand this turnaround and the dynamics that emerge within these programmes.

This overview not only demonstrates the attempt of the authors to combine theoretical rigour and 'requisite variety' but also the extension of topics to which principal–agent theory has been applied in science policy, though Shove sees some limitations in the approach. We could, however, question this conclusion, if we understand the relationship between programme directors and scientists as a relationship between a principal and multiple agents, and we use existing knowledge from the approach that underlines that such a constellation is only fruitful for the principal if the agents are in a competitive position. If, as is the case here, these agents play a co-operative game, it is not surprising that unique dynamics set in, and these agents use the opportunity structure of a research programme for their own purposes.

The articles in this special issue not only present the variety of topics that can be dealt with in terms of principal–agent theory, but they furnish also important insights such as:

- the importance of contextual features of the system that mitigate principal–agent relationships;

- the crucial role of funding agencies among these contextual features;
- the role of networks as a new way to organise principal–agent relations;
- the pertinence of the inclusion of users as the 'fourth actor';
- the conceptualisation of science policy in terms of configuration and not isolated bivariate principal–agent relationships; and
- the usefulness of institutional design in the organisation of science policy.

We, therefore, are convinced of the usefulness of the approach and would invite other scholars to embark on this quest for a more analytical and rigorous tool in observing science policy.

## References

Bates, R H, A Greif *et al* (1998), *Analytic Narratives.* (Princeton University Press, Princeton NJ).

Ben-David, J (1971), *The Scientist's Role in Society. A comparative study.* (Prentice-Hall, Inc, New Jersey).

Ben-David, J (1991), *Scientific Growth. Essays on the Social Organization and Ethos of Science* (University of California Press, Oxford).

Ben-David, J (1991 [1977]), "The central planning of science", in G Freudenthal (editor), *Joseph Ben-David. Scientific Growth* (University of California Press, Berkeley).

Bourdieu, P (1975), "The specificity of the scientific field and the social conditions of the progress of reason", *Social Science Information*, 14(6), pages 19–47.

Bourdieu, P (2001), *Science de la science et réflexivité* (Editions Raisons d'Agir, Paris).

Braun, D (1993), "Who governs intermediary agencies? Principal–agent relations in research policy-making", *Journal of Public Policy*, 13(2), pages 135–162.

Caswill, C (1998), "Social science policy: challenges, interactions, principals and agents", *Science and Public Policy*, 25(5), October, pages 286–296.

Coleman, J S (1990), *Foundations of Social Theory* (Belknap Press of Harvard University Press, Cambridge MA, London).

Daehle, van der W (1979), *Geplante Forschung.* (Suhrkamp TB, Frankfurt am Main).

Elgie, R (2002), "The politics of the European Central Bank: principal–agent theory and the democratic deficit", *Journal of European Public Policy*, 9(2), pages 186–200.

Ferejohn, J (1993), "Structure and ideology: change in Parliament in early Stuart England", in J Goldstein and R Keohane (editors), *Ideas and Foreign Policy. Belief system, Institutions and Political Change* (Cornell University Press, Ithaca).

Foray, D (2000), *L'économie de la connaissance* (La Découverte, Paris).

Ghiselin, M T (1987), "The economics of scientific discovery", in G Radnitzyk and P Bernholz (editors), *Economic Imperialism. The Economic Approach Applied Outside the Field of Economics.* (Paragon House Publishers, New York).

Gieryn, T F (1995), "Boundaries of science", in T J Pinch, *Handbook of Science and Technology Studies* (Sage, Thousand Oaks).

Guston, David H (1996), "Principal–agent theory and the structure of science policy", *Science and Public Policy*, 23(4), August, pages 229–240.

Guston, D H (2000), *Between Politics and Science* (Cambridge University Press, New York, Cambridge).

Huber, J, and A Lupia (2001), "Cabinet instability and delegation in parliamentary democracies", *American Journal of Political Science*, 45(1), pages 18–33.

Kiewiet, D R, and M D McCubbins (1991), *The Logic of Delegation. Congressional Parties and the Appropriation Process* (University of Chicago Press, Chicago).

Knorr, K D, R Krohn *et al* (1981), *The Social Process of Scientific*

*Investigation* (D Reidel Publishing Company, Dordrecht).

Knorr-Cetina, K D(1981), *The Manufacture of Knowledge. An Essay on the Constructivist and Contextual Nature of Science* (Pergamon Press, Oxford).

Krohn, W, and G Küppers (1987), "Die Selbstorganisation der Wissenschaft", Science studies report no 33, Universität Bielefeld.

Latour, B (1987), *Science in Action. How to follow scientists and engineers through society* (Harvard University Press, Cambridge MA).

Latour, B, and S Woolgar (1979), *Laboratory Life. The social construction of scientific facts* (Sage, London).

Luhmann, N (1990), *Die Wissenschaft der Gesellschaft* (Campus, Frankfurt am Main).

Majone, G (2001a), "Nonmajoritarian institutions and the limits of democratic governance: a political transaction-cost approach", *Journal of Institutional and Theoretical Economics*, 157, pages 57–78.

Majone, G (2001b), "Two logics of delegation. Agency and fiduciary relations in EU governance", *European Union Politics*, 2(1), pages 103–121.

Mayntz, R (1991), "Scientific research and political intervention — the structural development of publicly financed research in the Federal Republic of Germany", in F Orsi Battaglini and F Roversi Monaco (editors), *The University within the Research System — an International Comparison* (Nomos, Baden-Baden).

McCubbins, M D (1984), "Congressional oversight overlooked: police patrols versus fire alarms", *American Journal of Political Science*, 28, pages 165–179.

Merton, R K (1970 [1938]), *Science, Technology and Society in Seventeenth-Century England* (Harper and Row, New York).

Miller, G J (1992), *Managerial Dilemmas. The political economy of hierarchy* (Cambridge University Press, Cambridge).

Moe, T M (1984), "The new economics of organization", *American Journal of Political Science,* 28(4), pages 739–777.

Mukerji, C (1989), *A Fragile Power. Scientists and the State* (Princeton University Press, Princeton NJ).

Polanyi, M (1951), *The Logic of Liberty.*(Routledge and Kegan Paul, London).

Polanyi, M (1962), "The republic of science", *Minerva*, 1, pages 54–73.

Rip, A (1994), "The republic of science in the 1990s", *Higher Education Studies*, 28, pages 3–32.

Rip, Arie, and Barend J R van der Meulen (1996), "The post-modern research system", *Science and Public Policy*, 23(6), December, pages 343–352.

Ross, S A (1973), "The economic theory of agency: the principal's problem", *American Economic Review,* 12, pages 134–139.

Strom, K (2000), "Delegation and accountability in parliamentary democracies", *European Journal of Political Research*, 37(3), pages 261–289.

Tullock, G (1966), *The Organization of Inquiry* (Duke University Press, Durham NC).

Van der Meulen, B J R (1998), "Science policies as principal–agent games: institutionalization and path-dependency in the relation between government and science", *Research Policy*, 27, pages 397–414.

Weingart, P (1997), "From 'Finalization' to 'Mode 2': old wine in new bottles", *Social Science Information*, 36(4), pages 591–613.

Weingast, B E (1984), "The congressional–bureaucratic system: a principal–agent perspective (with application to the SEC)", *Public Choice*, 44, pages 147–191.

Williamson, O E (1975), *Markets and Hierarchies* (Free Press, New York).

Williamson, O E (1985), *The Economic Institutions of Capitalism. Firms, Markets, Relational Contracting* (The Free Press, New York).

# Agency Relationships and Governance Mechanisms in Service Delivery: A Theoretical Analysis[1]

## Debi P. Mishra[2]

## Abstract

The general topic of service quality has been widely studied in literature on marketing. Considered as a whole, researchers have focused on issues concerning the structure of service quality (e.g., SERVQUAL dimensions) and underlying psychological processes (e.g., role conflict, job stress) that impact delivery. While extant studies have added to our understanding of service quality, one notable gap in the literature concerns the lack of attention to *agency relationships* and *governance mechanisms* that affect delivery. For example, unless appropriate governance mechanisms or safeguards are in place, agents may *under-provide* or *over-provide* services, thereby adversely affecting quality. Given the widespread prevalence of agency relationships, the objective of this paper is to provide a focused discussion of agency problems and to specify how firms can deploy appropriate governance mechanisms to aid in the delivery of service quality.

## I. Introduction

The general area of service quality has received considerable attention in the marketing discipline. Up until the early and mid-1990's, the central focus of these studies was two-fold. First, researchers generated an impressive body of literature on the *structure* of the service quality construct by studying scale development, measurement, dimensionality, validity, and generalizability issues among others (Babukus and Boller, 1992; Babakus and Mangold, 1992). Second, a number of studies explored how individual level *psychological constructs* such as role conflict, role ambiguity, role stress, job burnout, and empowerment (Boulding et al., 1993; Cronin and Taylor, 1992) affected delivery. In recent years, the emphasis of the field has shifted somewhat with researchers focusing on behavioral and financial consequences of service quality.

While extant research has furthered our understanding of service delivery in a number of ways, one important gap in the literature remains unaddressed. Specifically, relatively less attention has been directed at understanding agency relationships that may serve as failure points by impacting quality negatively. An agency relationship is established whenever a principal hires an agent to do some work on the principal's behalf (Fama, 1980; Jensen and Meckling, 1983). The central problem in an agency relationship stems from information asymmetry or a situation where one party to the exchange such as the agent has more information than the principal (Bergen, Dutta, and Walker, 1992). In these situations, monitoring the agent becomes difficult and expensive and the agent may dilute quality.

Agency problems manifest themselves on a regular basis in service settings. For example, Mills (1990) observes that principals (i.e., patients and management) "are often unable to determine whether the tests and treatments of physicians are appropriate" (p. 35) which, in turn, affects service quality. A similar point is illustrated by a story involving Sears, a leading retailer in North America which provided unnecessary service to its customers (*Wall Street Journal*; October 2, 1992). Since auto-repair service is an experience good (Biehal, 1983; Nelson, 1974), customers (principals) who authorized all estimates prepared by Sears' mechanics (agents) did not know whether estimates were inflated or not. According to agency theory, Sears' agents (mechanics) engaged in *moral hazard* (i.e., opportunistic behavior) in order to earn high commissions because

[2] Ph.D., Associate Professor of Marketing, School of Management, State University of New York, USA.

principals (Sears' management and final customers) could not effectively monitor agents (mechanics). Consequently, service quality at Sears was adversely affected by an inefficient monitoring system based on output control (i.e., commissions). Sears subsequently shifted from "the incentive compensation system that paid employees solely on the basis of amount of repairs customers authorized" to "a program based on quality instead of quantity" (*Wall Street Journal*; June 23, 1992, p. B1).

Mechanisms such as compensation systems are governance modes that can ameliorate agency problems. The general resolution of the agency problem involves the deployment of both ex-ante and ex-post governance mechanisms (Bergen, Dutta, and Walker, 1992). Ex-ante governance mechanisms such as screening, training, etc. of service providers address the adverse selection problem while ex-post strategies such as appropriate compensation schemes can ameliorate the moral hazard problem in service delivery.

While the agency relationship between managers and service providers is readily apparent, there is another level of agency relationship involving the final customer as the principal and the company's brand (and managers) as the agent. Since services are characterized by *information asymmetry* or a situation where "buyers (unlike sellers) are not fully informed about product quality" (Rao and Bergen 1992, p. 413), management has an incentive to signal a firm's reputation (Bloom and Reve, 1990; Fombrun and Shanley, 1990) to the final buyer by using quality cues (e.g., *price, warranties, certification, investment in firm specific assets, price premiums*; Bloom and Reve, 1990; Klein and Leffler, 1981; Shapiro, 1983). By reducing information asymmetry (through signals of quality), service companies attempt to manage agency relationships with the final buyer.

If customers use signals such as brand names to choose a service company, there is no guarantee that promised service will actually be delivered unless managers govern their agency relationship with service providers. For example, service quality can be compromised if an individual agent decides to behave opportunistically and dilute quality. To sum it up, delivery of desired quality levels in the marketplace is contingent upon the deployment and use of governance mechanisms that can manage multiple agency relationships involving: 1) companies and service providers, and 2) companies and customers.

Despite the need to study agency problems little conceptual effort has been directed at researching agency relationships and governance mechanisms in service delivery. For a long time, researchers have called for studying internal processes concerning service delivery, but have lagged behind in developing appropriate conceptual approaches for studying such problems. Given this well articulated need for more studies relating to internal processes, why has the field been slow to answer such calls? One reason may be the lack of attention to appropriate theory for guiding research in this area.

In the light of the preceding discussion, the objective of this paper is to provide a detailed theoretical analysis of multiple agency relationships and governance mechanisms in service delivery. It is hoped that this study will close a major "gap" in our understanding by moving away from the extant "black box" attitude that neglects agency problems in service delivery.

This paper is organized as follows. First, I describe the nature of services and discuss how agency relationships manifest themselves at *two* levels (i.e., company-final customer, and company-service provider) in a company. This is followed by the delineation and description of a set of governance mechanisms that can ameliorate such agency problems. Finally, I comment upon the implications of this study for service marketers and describe the scope for further research.

## II. The Nature of Services and Levels of Agency Relationships

### Information asymmetry

There is near unanimous agreement among scholars in *marketing* (Lovelock, 1983; Zeithaml, 1981), *organization theory* (Bowen and Jones, 1986) and *operations/strategic management* (Nayyar, 1990, 1992, 1993), that *intangibility* is *the* critical goods-services distinction from which all other differences emerge" (Zeithaml, Parasuraman, and Berry, 1985; p. 33). Intangibility

refers to a situation where services "cannot be seen, felt, tasted, or touched in the same manner as goods" (Zeithaml, Parasuraman, and Berry, 1985, p. 33). Since services are performances which cannot be easily evaluated, one party to the transaction (the service provider) usually possesses more information than others (management and final customer). Owing to information asymmetry, customers cannot *a-priori* evaluate a company's service, while management lacks objective criteria to evaluate service providers. These evaluation problems are a direct consequence of intangibility or information asymmetry between service providers, agents, and other entities such as managers and customers (Bowen and Jones, 1986). Note that in the context of a performance evaluation problem, information asymmetry is also termed performance ambiguity. For the purpose of this paper we will use the terms information asymmetry and performance ambiguity interchangeably.

Jones (1990) notes that information asymmetry or performance ambiguity "is particularly prevalent when the goods or services being purchased are intrinsically complex, and their quality can only be *evaluated after* purchase" (p. 24). The presence of performance ambiguity in the "client-firm interface" (Mills and Turk 1986) leads to agency (principal-agent) problems because "one party (the principal) engages another party (the agent) to undertake actions on his behalf in situations of information asymmetry" (Clark and McGuiness, 1987; p. 8).

Information asymmetry gives rise to *two* principal-agent levels in a service encounter, i.e., between the company and the final customer, and between the company and service providers. Information asymmetry presents difficulties for customers to evaluate a service even after consumption (Siehl, Bowen and Pearson, 1992), thereby providing management with an incentive to reduce information asymmetry for gaining competitive advantage (Nayyar, 1990). Furthermore, information asymmetry also makes it difficult (i.e., costly) for management to monitor and control service providers (Anderson and Oliver, 1987; Eisenhardt, 1985). Specifically, management may have "limited direct control" over the "quality of service that is delivered", because "when employees are delivering intangible services", "they are essentially acting alone (Bowen and Schneider, 1988; p. 65).

To sum it up, by incorporating the concept of information asymmetry, researchers can explicitly focus on two levels of agency problems that are found in service companies.

## II (a). Agency Relationship between Company and Final Customers and Related Governance Mechanisms

As Bergen, Dutta, and Walker (1992) note, the "ultimate customer (principal) can be viewed as engaged in an agency relationship" with the company (agent) because services are performances (Holmstrom, 1985) which cannot be easily evaluated. Typically, sellers have more information than buyers (information asymmetry) about the true quality of the service. This information asymmetry can lead to "moral hazard", because the company may exert less than complete effort in providing the service, or it may overprovide the service. Though customers (principals) attempt to reduce this information asymmetry by relying on "word-of-mouth" communications (Biehal, 1983; Murray, 1991) or by "purchasing" cheap information on the agent (company) from institutions (e.g., surrogate customers: Solomon, 1986; consumer reports: Hill and Jones, 1992), management has an incentive to reduce information asymmetry for the final buyer by using signals.

There are three main reasons for firms to signal quality and reduce customers' adverse selection risks. *First*, service customers "seek risk-reducing" information because of the "intangible, ephemeral, and experiential nature of services" (Murray, 1991, p. 20) in order to make better choices (Stigler, 1961). In fact, as Fombrun and Shanley (1990) observe, "the more informational asymmetry and ambiguity characterize the interactions between management and stakeholders (customers), the more likely the latter are to search for information" (p. 235). Consequently, "service firms can develop competitive advantage by exploiting the potential buyer's incentives to lower information acquisition costs when buying services" (Nayyar, 1990; p. 513). Companies attempt to reduce information asymmetry by providing customers with surrogate barometers of quality (Akerlof, 1970). These surrogates may be considered as *signals* and defined as "marketer-controlled easy-to-acquire informational cues, extrinsic to the products themselves, that consumers use to form inferences about the quality or value of those products" (Bloom and Reve, 1990; p.

59). More technically, using an agency theory perspective, Cooper (1992) comments upon the function of signals as follows:

> *In many markets, one agent has private information that could help others in making their decisions. The uninformed agents would usually adjust their actions to suit their environment better if they could learn the private information before making choices. Because of this potential to change actions, sharing the private information could benefit the more informed agents or society as a whole. One method of disclosing private information is signaling (p. 431; emphases added).*

*Second,* by reducing information asymmetry, service companies prevent market failure and contribute to social good. For instance, some service companies may exploit information asymmetry to their advantage (by engaging in moral hazard) and supply low quality services. These firms may earn supernormal profits (because of low production costs), which provide no incentive to honest firms for staying in business. In the extreme case, "it is quite possible to have the bad driving out the not-so-bad driving out the medium driving out the not-so-good driving out the good in such a sequence of events that no market exits at all" (Akerlof, 1970; p. 490). This has been termed the "lemons" problem Akerlof (1970).

*Finally*, service firms which earn a good reputation by reducing performance ambiguity can successfully diversify into related services (Nayyar, 1990) by legitimately transferring their reputation to new services. According to Nayyar (1990), "a firm that diversifies into services that its existing customers may buy could create a competitive advantage, since it could potentially exploit the favorable attention in the information asymmetry distribution faced by potential buyers when they consider buying the new service offered by the firm" (p. 516).

### Governance mechanisms

The general strategy to solve customer's agency problems is called signaling. Signaling strategies which firms may use to reduce information asymmetry for the final buyer are of two types, i.e., *direct* and *indirect* (Nayyar, 1990). Direct quality signals assure the buyer of a minimal level of performance by reducing information asymmetry. The most widely mentioned signals are guarantees and certification. Guarantees shift the risk of purchase from the buyer to the seller and ensure some level of quality (Akerlof, 1970). According to Hill and Jones (1992), some services are inherently difficult for buyers to evaluate prior to purchase. The existence of information asymmetry presents customers "with a difficult agency problem" because "the consumer is vulnerable to opportunistic action on the part of management and the agency problem is solved by the *ex-ante* introduction of a warranty into the contracting scheme" (p. 139). The use of warranty as an information asymmetry reduction mechanism has also been suggested by Allen (1984), Grossman (1981), and Wiener (1985).

Though some authors (Bergen, Dutta, and Walker, 1992) suggest that the efficacy of guarantees is limited because of customers' proclivities to behave opportunistically (e.g., by falsifying a claim), service marketers (Hart, Schleisinger, and Maher, 1992) have stressed the power of unconditional service guarantees. More importantly, service guarantees offered to final buyers also act as a vehicle for communicating quality levels to employees. Furthermore, the presence of an unconditional guarantee like "customer satisfaction" provides management with objective criteria for monitoring boundary spanners (frequency with which guarantees are invoked), whose behavior is typically difficult to observe.

Service guarantees serve to reduce the "gap" (Zeithaml, Parasuraman, and Berry, 1988) between management and boundary spanners about quality perceptions. In other words, by offering guarantees, management attempts to solve not only the agency problem with final customers but also the agency problem with service providers. Service guarantees therefore differ from product guarantees which are directed solely at the final buyer and are attempts to solve only one level of agency problems (between management and the final buyer). In this vein, Nayyar's (1990) observation that "warranties covering services are impossible to administer since failure to perform a social interaction is generally indeterminable" (p. 514), ignores the potential of service guarantees to solve the agency problem between management and service providers.

Certification, which indicates the "attainment of levels of proficiency", also reduces quality uncertainty (Akerlof, 1970; p. 500). According to Akerlof (1970), "the high school diploma, the baccalaureate degree, the Ph.D., and even the Nobel Prize, to some extent, serve this function of certification" (p. 500). By prominently stressing the qualifications of their professors, universities seem to reduce performance ambiguity for freshmen.

Indirect quality signals serve to reduce information asymmetry for the final buyer by stressing a firm's reputation. Klein and Leffler (1981) suggest two such signals, i.e., *price premiums* and *firm-specific capital investments.* Firms can signal high quality by charging prices above the market price (i.e., charging a price premium). However, if these firms cheat on quality, a potential stream of future profits would be lost. According to Klein and Lefler (1981), "this price premium stream can be thought of as protection money paid by consumers to induce contract performance" (p. 624). Thus, when firms do not deliver the promised level of quality, customers may withdraw this deposit, causing the firm to go out of business. Price premiums indirectly reduce information asymmetry for buyers by promising quality. A firm charging price premiums have every incentive to maintain the quality of its services and reap future profits which are held hostage in view of possible quality dilution.

Firm specific capital investments yield only "small direct consumer services… relative to cost" (Klein and Leffler, 1981; p. 627). For instance, expensive advertising for services characterized by high levels of information asymmetry does not necessarily reveal relevant information (1974). As an illustration, hospitals' advertisements do not detail the surgical procedure for patients with heart problems. On the other hand, the purpose of expensive and "non-informative" advertising (Nelson, 1974) is to signal a company's reputation to the final buyer and to reduce information asymmetry.

Firms making company specific investments trade off "increased consumer service value with decreased salvage value" (Klein and Leffler, 1981; p. 627). According to these authors, "the expenditures on brand name capital assets are similar to collateral that a firm loses if it supplies output of less than anticipated quality (p. 627). Examples of firm specific capital investments are *logos and expensive signs, ornate settings like expensive carpets and upholstery which yield no direct service, human entrepreneurial skills and idiosyncratic knowledge, expensive advertising, and celebrity advertising* (Klein and Leffler, 1981; Rubin, 1990).

Interestingly, it has been recognized that the use of ornate settings or "elaborate service-scapes" (Bitner, 1992) is an attempt to make the service ambience physiologically pleasant for the customer (Bitner, 1992). However, agency theory suggests that ornate settings in hospitals are signals of reputation which management uses to solve the agency problems with patients. In other words, patients realize that hospitals have sunk a lot of money into these expensive investments (e.g., ornate settings). Consequently, firms cannot possible cheat on quality.

In a study on the relationship between reputation effects and price premiums, Rao and Bergen (1992) found out that reputable sellers could not command price premiums. One possible explanation for this finding is that these sellers did not make commensurate investments in firm specific assets. As Klein and Leffler (1981) note, firms may command price premiums only when commensurate investments have been made in firm specific assets. In other words, buyers may not pay high prices to "seemingly" reputable agents who do not make collateral investments, fearing a "rip-off" (Dejong, Forsythe, and Lundholm, 1985).

Though a number of signaling strategies have been suggested in the literature (e.g., warranties, certification, investments), existing theory does not comment upon the relative importance of these signals for solving agency problems. Echoing this point, Rao and Bergen (1992) note that "future research will be required to suggest which of these many devices is most appropriate for a given situation" (p. 421). There is some discussion in literature on strategic management (Nayyar, 1990) that firms may focus more on indirect signals of quality (e.g., reputations) than on warranties and certification. As Nayyar (1990) notes, "certification, too, is so widely prevalent as to make it of no consequence in consumer choice behavior" (p. 514). Perhaps companies realize that reputation is an idiosyncratic asset (Rashid, 1988) which cannot be easily duplicated by competitors.

## II (b). Agency relationship Between Company and Service Providers and Related Governance Mechanisms

Service is finally delivered to customers by boundary spanning employees (Aldrich and Herker, 1977). When performance ambiguity is high, management (principal) cannot completely and costlessly monitor the actions of service providers (agents). Due to incomplete monitoring, service providers may engage in moral hazard (opportunistic behavior) and oversupply or undersupply the service, thereby adversely affecting service quality. Moral hazard is a typical problem when services are high in credence properties (e.g., medical care and education; Darby and Karni, 1973). For instance, Swedlow et al. (1992) note that "MRI (Magnetic Resonance Imaging) scans (were) medically inappropriate 38% more often when ordered by self-referring physicians, suggesting increased rates of use in this group" (p. 1506). In a similar vein, Gomez-Mezia and Balkan (1992) observe that "in a university setting, principals face a classical agency problem with respect to faculty" and that "information asymmetries between faculty and administrators (principals) create steep agency costs for the latter if they attempt to directly monitor faculty behavior" (p. 923). Furthermore, most professors in universities have a lot of freedom in designing courses and conducting research. There is a possibility that a professor may engage in moral hazard by putting in less effort into teaching and research than into consulting. University administrators face therefore the classic agency problem of preventing "faculty members (agents) from taking advantage of their privileged and nonprogrammable position" (Gomez-Mezia and Balkan, 1992; p. 924).

### Governance mechanisms

Management solves the agency problem with boundary spanning employees by using various types of control mechanisms (i.e., output and behavior controls: Eisenhardt, 1985; Ouchi, 1980). One approach to managing such information asymmetry entails the resolution of *adverse selection* or *hidden action* problems (Bergen, Dutta, and Walker, 1992).

According to the adverse selection model, managerial strategies depend on the extent to which agent's actions can be costlessly observed. In general, when management can observe the behavior of boundary spanners easily, behavior control strategies (e.g., hourly pay systems) are suggested. This method of control is suitable for services characterized by low information asymmetry (e.g., grocery stores, where a sales clerk's actions are routinized). On the other hand, agency theory recommends the use of output control (e.g., commissions) when employee behavior cannot be costlessly observed. However, for highly intangible services (e.g., medical care and education), a commission system may be inappropriate because it places little emphasis on customer satisfaction (Anderson and Oliver, 1987). Accordingly, management uses complex compensation systems for aligning the interests of service providers with those of the company. For instance, Gomez-Mezia and Balkan (1992) note that university administrators often tie a professor's compensation to the number of quality journal articles he or she publishes. Likewise, hospitals may link the bonuses of physicians to "patient satisfaction" scores (Dranove and White, 1987).

It should be noted that applications of agency theory in marketing have concentrated rather narrowly on "salesperson's compensation" issues (John and Weitz, 1989; Lal and Staelin, 1986; Oliver and Weitz, 1991). Practically no attention has been directed at *ex-ante* strategies which management can use (e.g., rigorous screening of employees) to prevent *ex-post* contractual problems (e.g., shirking by service providers). In agency theory parlance, *hidden information or adverse selection* strategies (e.g., screening, socialization, and training: Bergen, Dutta, and Walker, 1992) have not been researched. For intangible services, management often uses rigorous screening procedures to ensure that service providers' subsequent performances are congruent with company objectives. For example, university professors at the entry level are selected through a rigorous process which involves several rounds of screening (preliminary screening, initial interview, campus visits and presentations and careful consideration of reference letters). By following this extended search procedure, universities try to discover as much "hidden information" as they can on a candidate prior to his or her selection. Another example of how service companies rigorously select and train service providers is provided by an illustration in the Wall Street Journal (January 25, 1993):

*Kaiser Permanente, a health maintenance organization… is often cited as a model health plan for managed competition. It recruits doctors through an evaluation process that includes a rigorous review of training and credentials and hires them as probationary em- ployees for three years. At the end of that period, doctors are voted in as full-fledged mem- bers by their peers, based on advanced training, perceptions of competence as such factors as rapport with patients and staff (p. A12).*

Note that although Gomez-Mezia and Balkan (1992) correctly view the relationship be- tween administrators and professors as an agency problem, they focus only on "compensation". In other words, though "hidden action" or moral hazard problems have been addressed in the litera- ture, "hidden information" issues have received less attention. In the context of services, it is im- perative to research both "hidden action" and "hidden information" models.

The findings from agency studies on salesforce compensation plans (Oliver and Weitz, 1991) are relevant to service organizations. However, some clarifications are in order. In the sales- force literature, the exogenous concept of environmental uncertainty determines subsequent com- pensation plans (e.g., salary or commissions) for boundary spanners. Environmental uncertainty is often operationalized as uncertainty in the relationship between effort expended and sales (results) (Oliver and Weitz, 1991). John and Weitz (1989) measure uncertainty of "product sales" as an indicator of environmental uncertainty. The focus on "sales" inevitably excludes any consideration of services as salespeople may overprovide services in order to earn high commissions, thereby affecting customer satisfaction and service quality. The Sears situation discussed earlier vividly illustrates the adverse effect of "commission" systems on service quality. According to agency theory, the use of output control systems for salespeople is less effective when environmental un- certainty is high, because agents are assumed to be "risk averse". Being risk averse, agents facing a highly uncertain environment will not opt for a commission system. In other words, they are better off with some assured compensation (e.g., salary). In sum, agency theory predicts that when envi- ronmental uncertainty is high, salary based systems are effective because of the "risk averse" na- ture of agents. On the other hand, when environmental uncertainty is low, salary is also the domi- nant compensation mode because an agent's behavior can easily be observed (by management).

The findings from compensation schemes for salespeople is directly applicable to service settings although it is important to recognize that "risk aversion" plays no part in determining compensation of service providers. For instance, when information asymmetry is low, manage- ment can easily observe an agent's behavior and a 'salary' system is recommended. When infor- mation asymmetry is high, output based systems are inadequate because agents may engage in 'moral hazard' and overprovide or underprovide services to final customers. Hence, when infor- mation asymmetry is high, service providers may be compensated with 'salary' not because they are 'risk averse', but because they may provide poor service to the final customer. Furthermore, when information asymmetry is high, a complex compensation system for agents which incorpo- rates notions of service quality and customer satisfaction may also be used. In any case, output control systems are clearly inappropriate in a service setting because they can compromise service quality – a point recognized by Anderson and Oliver (1987).

In sum, managerial strategies for solving the agency problem with service providers should include elements from both the "hidden action or moral hazard" and the "hidden informa- tion or adverse selection" models. Perhaps less attention has been paid to "hidden information" strategies because they have traditionally been considered outside the domain of "agency" theory (Eisenhardt, 1985). As such, the incorporation of "hidden information" models into agency theory in recent studies (Bergen, Dutta, and Walker, 1992) is a welcome trend.

### Interdependencies between agency relationships

Although agency relationships in a service organization exist at two levels, they are inter- dependent. Management solves the agency relationship with final customers by using signals of reputation. Reputation is defined as "a set of attributes ascribed to a firm, inferred from the firm's past actions" (Weigelt and Camerer, 1988; p. 443). Service firms assure customers of quality by stressing reputations. Klein and Leffler (1981) argue that investments in firm specific assets are

essentially reputation building activities which serve to assure buyers about quality. In other words, by compromising on quality, these firms risk the appropriation of future quasi-rents (Klein, Crawford, and Alchain, 1978). Rashid (1988) articulates this point well by noting that "when significant amounts of money are invested, the businessman tells that he plans to stay for some time to come… in the long run the only way to stay is by pleasing customers… this requires providing them with the goods they really want… *this long-term dependence of producers on customers is perhaps the most effective guarantee of quality*" (p. 248; emphasis added).

According to Camerer and Vepsalainen (1988), managers have an incentive to maintain the reputation of their firms. The owner-manager has an incentive to maintain his firm's reputation so as to increase its salvage value. On the other hand, in the case of firms where ownership and control are separated, "managers are continuously 'selling' the firm to new owners through capital markets… managers who erode the firm's reputation are depreciating an intangible asset and are vulnerable to market discipline like takeover attempts" (p. 118).

Maintaining a firm's reputation solves the agency problem between management and the final customer because reputation is essentially an information asymmetry reduction strategy. On the other hand, service is actually delivered by a distant boundary spanner who may act in his or her self interest and compromise on quality. Management is therefore faced with a problem of safeguarding its reputation because it is involved in a second agency relationship with the service provider. In this sense, the two agency levels appear to be inter-related. In other words, the greater the use of reputation by a management for reducing performance ambiguity for the final customer is, the greater the need to monitor service providers appears to be. Klein and Leffler (1981) discuss this problem in the context of a franchisor-franchise relationship:

> *The existence of independent competitive retailers that do not have any ownership stake in this firm specific asset and yet can significantly influence the quality of the final product supplied to consumers creates a severe quality-cheating problem for the manufacturer. Manufacturers may protect their trademarks by imposing constraints on the retailer competitive process including entry restrictions, exclusive territorial grants, minimum resale price maintenance, and advertising restrictions that will assure quality by creating a sufficiently valuable premium stream for retailers" (Klein and Leffler, 1981; p. 633).*

In a similar vein, Brickley and Dark (1987) argue that "a major problem facing companies with valuable names is controlling the action of agents throughout the organization to assure the continued value of that trademark" (p. 403). An example of how reputation effects can be compromised in a service setting because of agency problems between management and service providers is illustrated by Dejong, Forsythe, and Lundholm (1985). These authors explicitly model reputation effects in studying a principal-agent problem in the stock market. The findings of this study indicate that "while there is evidence of reputation effects in these markets seemingly reputable agents are often able to use opportunities for false advertising to their advantage and 'rip-off' principals" (p. 809). To sum it up, reputation effects alone do not guarantee quality because "the presence of moral hazard does indeed lead to the provision of nonoptimal levels of services in an agency relationship (between management and the service provider) (Dejong, Forsythe, and Lundholm, 1985, p. 819).

The notion of multiple agency relationships, their interrelated nature, and the impact of governance mechanisms on service delivery are depicted in Figure 1. Panel A of the figure depicts the simplest possible agency relationship between a single principal and a single agent. For example, when a patient (principal) obtains service from an independent physician (agent), such a relationship is established. From the principal's standpoint, the two main agency problems that need to be governed are as follows: i) adverse selection, and ii) moral hazard. Efforts undertaken by the patient to pre-qualify an independent physician such as screening, word-of-mouth referrals from other patients, etc., constitute governance mechanisms for resolving the adverse selection problem. Moral hazard or hidden action problems may come up after the patient has begun visiting the physician on a regular basis. Governance mechanisms that ameliorate moral hazard may have to do with the length of the doctor-patient relationship and the building up of trust. As such, doctor-patient relationships are often sticky because principals having resolved adverse selection and

moral hazard problems do not wish to grapple with additional uncertainty by switching to a new physician. In this setting, using compensation as a governance mechanism appears less relevant given institutionalized compensation practices.

Next, consider panels B and C of Figure 1 that depict multiple agency relationships. The key idea is that multiple agency relationships are i) interdependent, and ii) need to governed simultaneously in order to yield optimal quality outcomes. In Panel B, consider a situation where McDonald's as the franchiser (principal) deals with the individual franchisee (agent). McDonalds faces both adverse selection and moral hazard problems because individual franchisees that are far removed from headquarters may dilute delivered quality. These problems are managed in several ways. In addition to instituting governance mechanisms like pre-qualification and training of franchisees, McDonald's also employs District Sales Managers (DSM's) who undertake field visits to monitor individual franchisees. Hence, the DSM is a principal in his or her relationship with the franchisee and is an agent of McDonalds. To deliver optimal levels of quality, McDonald's has to craft appropriate governance safeguards at the level of the DSM also, e.g., through pre-qualification and compensation mechanisms. However, in governing the relationship with DSM's (Level 1), agency relationships at another level (Level 2) may also be affected. For example, if the DSM is compensated on the basis of sales, he or she, in turn, may impose additional burden on franchisees by imposing unattainable sales goals on them. As a consequence, individual franchisees can get de-motivated and may switch to competition.
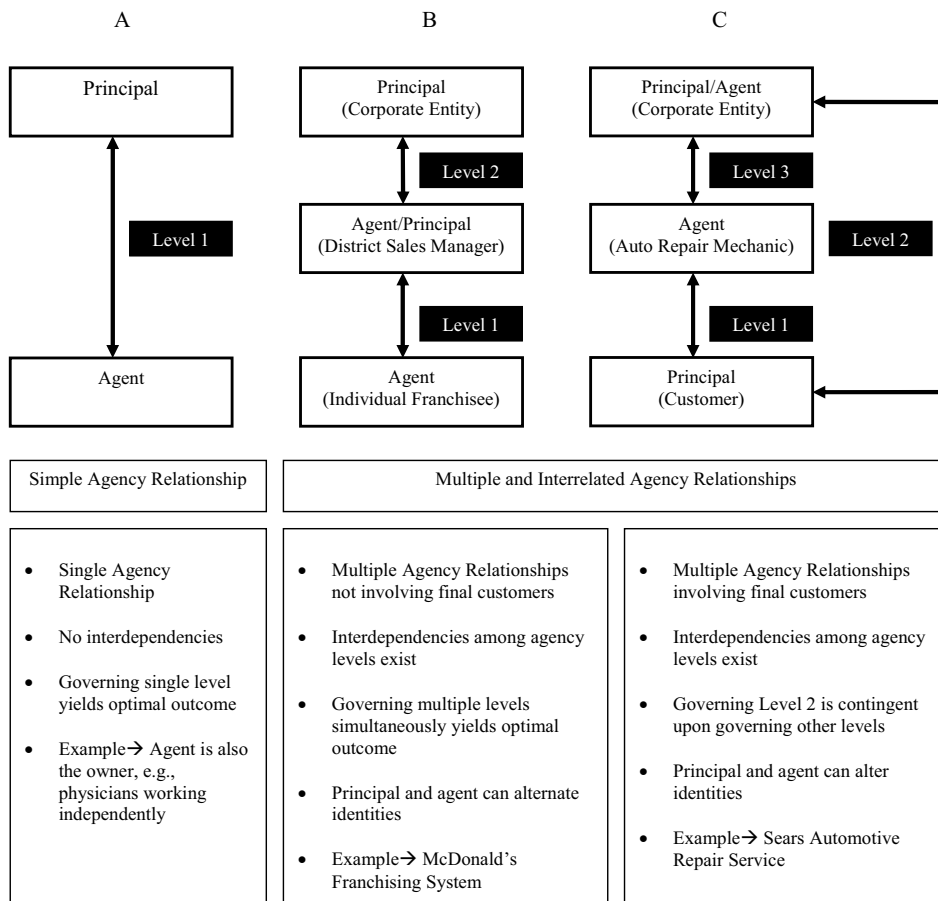


Fig. 1. Simple and Interdependent Agency Relationships in Service Delivery

Panel C depicts a situation described earlier. Typically, automotive repair chains such as Firestone, Midas, Goodyear, Sears, etc. invest considerable resources in promoting their brand. Brand image, in turn, acts as a quality signal and governs customers' adverse selection concerns. In general, other things being equal, customers are more likely to choose a brand that is widely advertised and has a national presence. However, governing this level of agency relationship (Level 2) is contingent upon how other agency relationships are managed. In the simplest case, promise quality at level 2 may be diluted if management has not solved the agency relationship with mechanics or managers who are directly involved in providing the service (Level 3). In sum, optimal delivery of service quality entails i) the crafting of appropriate governance mechanisms, and ii) the resolution of interdependency problems.

## III. Implications And Scope For Further Research

It was argued at the outset that agency problems in service delivery have not received any systematic attention in marketing. This gap in the services literature is surprising because most service arrangements in today's society are agency relationships. Specifically, in keeping with the transition from an agrarian economy to industrial one, role specialization has engendered the modern agency problem. To address this shortcoming, we used agency theory to study multiple levels of agency relationships and specified appropriate governance mechanisms. This study has implications for managers and researchers.

### *Managerial implications*

Managers should realize that service quality involves more than a smile or a handshake. Quality is the result of a process which starts within the organization. Service organizations are characterized by the presence of agency relationships where parties often have divergent interests. The successful resolution of agency problems at different levels within an organization is the *sine-qua-non* for achieving quality.

Despite the obvious importance of understanding agency relationships, prescriptions from extant research can be summarized in the following sentence: *These are the dimensions of quality, now manage your internal activities in accordance with these dimensions to that you achieve that elusive mantra for profitability – quality.* In sum, no systematic understanding exists in the literature on how managers should manage agency relationships in which they are involved. For instance, the Sears example discussed earlier is a classic case of mismanagement of an agency relationship where two principals (Sears' management and the final customer) and an agent (Sears' mechanics) were locked in an inefficient arrangement that compromised quality.

By viewing monitoring and control systems through the powerful lens of agency theory, managers can be better equipped to solve these problems. First, by classifying their service along the dimensions of information asymmetry, managers will be sensitized to the relative importance of monitoring problems. Second, the recognition of agency problems may help management to transmit signals to final customers. For example, management may realize that it cannot charge premium prices for its service without making appropriate investments in improving the ambience of the service setting. Quality conscious customers will perceive these specific investments as "collaterals" for price premiums. Advertising strategies which stress the company's reputation may also send powerful quality signals to customers. Furthermore, if guarantees are widespread, management can gain competitive advantage by introducing unconditional guarantees. Third, appreciation of agency problems may motivate management to design appropriate compensation schemes for service providers. For example, a dentist's compensation may be tied to patient satisfaction scores. In this way, management can dovetail its monitoring of service providers with its objective of providing superior quality.

### *Scope for further research*

An obvious avenue for future research is to formulate empirically testable hypotheses in the context of extant theory. Researchers may also incorporate "organizational culture" concepts (Deshpande and Webster, 1989; Deshpande, Farley, and Webster, 1993) to better understand

agency problems. Attention to culture issues is desirable because shared ideas and beliefs can minimize the divergent interests of parties involved in an agency relationship. For example, researchers can empirically determine the relative importance of agency strategies and organizational culture in delivering service quality. The results of such an empirical study may directly impinge upon the debate in *marketing* and *organization behavior* about the narrow focus of agency theory on human opportunism. There is some empirical evidence in literature on marketing (Heide and John, 1992) which suggests that innate human values like trust and norms may act as safeguards against human opportunism. In fact, in some organizations, culture may be the most important determinant of service quality.

Finally, the conceptual framework offered in this study can be extended to better understand various "forms" of service organizations. For instance, the health care industry has different types of organizational forms, e.g., *fee-for-service, autonomous physicians, health maintenance organizations, referral networks.* In the fee-for-service form, the agency arrangement between a doctor and his or her patient is rendered efficient because principals develop close relationships with agents. This relationship assures the doctor of continued loyalty and future business. On the other hand, health maintenance organizations manage agency relationships with doctors through various compensation schemes and socialization strategies. In this vein, we may note that the relationship between a customer and a "surrogate customer" (e.g., expert services like stock-analysts) is not a principal-agent relationship. By employing a surrogate customer, the buyer attempts to obtain cheap information on *the* agent, i.e., on the company where he or she desires to buy the service. To sum it up, many forms of service delivery in today's business environment can be studied by using principles derived from agency theory.

# References

1. Akerlof, George A., (1970), "The Market for Lemons: Quality Under Uncertainty and the Market Mechanism," *Quarterly Journal of Economics*, 84 (August), 488-500.
2. Aldrich, H.E., and D. Herker (1977), "Boundary Spanning Roles and Organizational Structure," *Academy of Management Review*, 2, 217-230.
3. Allen, Franklin (1984), "Reputation and Product Quality," *Rand Journal of Economics*, 15 (3), Autumn, 311-327.
4. Anderson, Erin, and Richard Oliver (1987), "Perspectives on Behavior-Based Versus Outcome-Based Control Systems," *Journal of Marketing*, 51 (October), 76-88.
5. Babakus, Emin, and Gregory W. Boller (1992), "An Empirical Assessment of the SERVQUAL Scale," *Journal of Business Research*, 24, 253-268.
6. Babakus, Emin, and W. Glynn Mangold (1992), "Adapting the SERVQUAL Scale to Hospital Services: An Empirical Investigation," *Health Services Research*, 26 (February), 6, 767-786.
7. Bergen, Mark E., Shantanu Datta, and Orville C. Walker Jr., (1992), "Agency Relationships in Marketing: A Review of the Implications and Applications of Agency Related Theories," *Journal of Marketing*, 56 (3), July, 1-24.
8. Biehal, Gabriel J. (1983), "Consumers' Prior Experiences and Perceptions in Auto Repair Choice," *Journal of Marketing*, 47 (Summer), 82-91.
9. Bitner, Mary Jo (1992), "Servicescapes: The Impact of Physical Surroundings on Customers and Employees," *Journal of Marketing*, 56 (April), 57-71.
10. Bloom, Paul N., and Torger Reve (1990), "Transmitting Signals to Consumers for Competitive Advantage," *Business Horizons*, July-August, 58-66.
11. Boulding, William, Ajay Kalra, Richard Staelin, and Valarie A. Zeithaml (1993), "A Dynamic Process Model of Service Quality: From Expectations to Behavioral Intentions," *Journal of Marketing Research*, 30 (February), 7-27.
12. Bowen, David E., and Benjamin Schneider (1988), "Services Marketing and Management: Implications for Organization Behavior," in *Research in Organization Behavior*, 10, 43-80.
13. Bowen, David E., and Gareth R. Jones (1986), Transaction-Cost Analysis of Service Organization-Customer Exchange," *Academy of Management Review*, 11 (2), 428-441.

14. Brickley, James A., and Frederick H. Dark (1987), "The Choice of Organizational Form: The Case of Franchising," *Journal of Financial Economics*, 18, 401-420.

15. Camerer, Colin, and Ari Vepsalainen (1988), "The Economic Efficiency of Corporate Culture," *Strategic Management Journal*, 9, 115-126.

16. Clark, Roger and Tony McGuiness (1987), "Introduction," in *The Economics of the Firm*, Roger Clarke and Tony McGuiness (Eds.), 1-17, New York: Basil Blackwell.

17. Cronin, Joseph J. and Steven A. Taylor (1992), "Measuring Service Quality: A Re-examination and Extension," *Journal of Marketing*, 56 (July), 55-68.

18. Darby, M.R. and E. Karni (1973), "Free Competition and the Optimal Amount of Fraud," *Journal of Law and Economics*, 16 (April), 67-86.

19. Dejong, Douglas V., Robert Forsythe and Russell J. Lundholn (1985), "Ripoffs, Lemons, and Reputation Formation in Agency Relationships: A Laboratory Market Study," *Journal of Finance*, XL (3), July, 809-823.

20. Deshpande, Rohit, John U. Farley and Frederick E. Webster (1993), "Corporate Culture, Customer Orientation, and Innovativeness in Japanese Firms: A Quadrad Analysis," *Journal of Marketing*, January.

21. Deshpande, Rohit and Frederick E. Webster (1989), "Organization Culture and Marketing: Defining the Research Agenda," *Journal of Marketing*, 53 (January), 3-15.

22. Dranove, David and William D. White (1987), "Agency and the Organization of Health Care Delivery," *Inquiry*, 24 (Winter), 405-415.

23. Eisenhardt, Kathleen M. (1985), "Control: Organizational and Economic Approaches," *Management Science*, 31, 134-149.

24. Fama, Eugene F. (1980), "Agency Problems and The Theory of the Firm," *Journal of Political Economy*, 88 (21), 288-307.

25. Fombrun, Charles and Mark Shenley (1990), "What's in a Name? Reputation Building and Corporate Strategy," *Academy of Management Journal*, 33 (2), June, 233-258.

26. Gomez-Mejia, Luis R. and David B. Balkan (1992), "Determinants of Faculty Pay: An Agency Theory Perspective," *Academy of Management Journal*, 35 (5), December, 921-955.

27. Grossman, Sanford (1981), "The Informational Role of Warranties and Private Disclosure about Product Quality," *Journal of Law and Economics*, 24 (December), 461-483.

28. Hart, Christopher W.L., Leonard L. Schleisinger and Dan Maher (1992), "Guarantees Come to Professional Service Firms," *Sloan Management Review*, Spring, 19-27.

29. Heide, Jan B. and George John (1992), "Do Norms Really Matter in Marketing Relationships?" *Journal of Marketing*, 56 (April), 32-44.

30. Hill, Charles W. and Thomas M. Jones (1992), "Stakeholder-Agency Theory," *Journal of Management Studies*, 29 (2), March, 131-154.

31. Holmstrom, B. (1985), "The Provision of Services in a Market Economy," in R.P. Inman (Ed.), *Managing the Service Economy: Prospects and Problems:* 183-213, Cambridge, U.K.: Cambridge University Press.

32. Jensen, Michael and W. Meckling (1983), "Theory of the Firm: Managerial Behavior, Agency Costs, and Capital Structure," *Journal of Financial Economics*, 3 (October), 305-360.

33. John, George and Barton A. Weitz (1989), "An Empirical Investigation of Factors Related to the Use of Salary Versus Incentive Compensation," *Journal of Marketing Research*, 26 (February), 1-14.

34. Jones, Gareth R. (1990), "Governing Customer-Service Organizational Exchange," *Journal of Business Research*, 20, 23-29.

35. Klein, Benjamin and Keith B. Leffler (1981), "The Role of Market Forces in Assuring Contractual Performance," *Journal of Political Economy*, 89 (4), 615-641.

36. Klein, Benjamin, R. Crawford and Armen Alchian (1978), "Vertical Integration, Appropriable Rents, and the Competitive Contracting Process," *Journal of Law and Economics*, 21: 297-326

37. Lal, Rajiv and Richard Staelin (1986), "Salesforce Compensation Plans with Asymmetric Information," *Marketing Science*, 5 (Summer), 179-198.

38. Lovelock, Christopher H. (1984), *Services Marketing*, Engelwood Cliffs, NJ: Prentice-Hall.

39. Mills, Peter K. (1990), "On the Quality of Services in Encounters: An Agency Theory Perspective," *Journal of Business Research*, 20, 31-41.
40. Mills, Peter K. and T. Turk (1986), "A Preliminary Investigation into the Influence of Customer-Firm Interface on Information Processing and Task Activities in Service Organizations," *Journal of Management*, 12 (1), 91-104.
41. Murray, Keith B. (1991), "A Test of Services Marketing Theory: Consumer Information Acquisition Activities," *Journal of Marketing*, 55, (January), 10-25.
42. Nayyar, Praveen R. (1990), "Information Asymmetries: A Source of Competitive Advantage for Diversified Service Firms," *Strategic Management Journal*, 11, 513-519.
43. Nayyar, Praveen R. (1992), "Performance Effects of Three Foci in Service Firms," *Academy of Management Journal*, 35 (5), 985-1009.
44. Nayyar, Praveen R. (1993), "Performance Effects of Information Asymmetry and Economies of Scope in Diversified Service Firms," *Academy of Management Journal*, 36 (1), 28-57.
45. Nelson, Phillip (1970), "Information and Consumer Behavior," *Journal of Political Economy*, 81 (4), July-August, 729-754.
46. Nelson, Phillip (1974), "Advertising as Information," *Journal of Political Economy*, 81 (4), July-August, 729-754.
47. Oliver, Richard L. and Barton A. Weitz (1991), "The Effects of Risk Preference, Uncertainty, and Incentive Compensation on Salesforce Motivation," *Report Number 91-104*, Cambridge, MA: Marketing Science Institute (February).
48. Ouchi, William G. (1980), "Markets, Bureaucracies, and Clans," *Administrative Science Quarterly*, 25 (March), 129-141.
49. Rao, Akshay R. and Mark E. Bergen (1992), "Price Premiums as a Consequence of Buyers' Lack of Information," *Journal of Consumer Research*, 19, 3 (December), 412-423.
50. Rashid, Salim (1988), "Quality in Contestable Markets: A Historical Problem?" *Quarterly Journal of Economics*, February, 246-249.
51. Rubin, Paul (1990), *Managing Business Transactions*, New York: John Wiley.
52. Shapiro, Carl (1983), "Premium for High Quality Products as Returns to Reputations," *Quarterly Journal of Economics*, 98 (November), 659-689.
53. Siehl, Caren, David E. Bowen and Christine M. Pearson (1992), "Service Encounters as Rites of Integration: An Information Processing Model," *Organization Science*, 3 (4), November, 537-555.
54. Solomon, Michael R. (1986), "The Missing Link: Surrogate Customers in the Marketing Chain," *Journal of Marketing*, 50 (October), 208-218.
55. Stigler, George J. (1961), "The Economics of Information," *Journal of Political Economy*, 69, 104-122.
56. Swedlow, Alex, Gregory Johnson, Neil Smithline and Arnold Milstein (1992), "Increased Costs and Rates of Use in the California Worker's Compensation System as a Result of Self-Referral by Physicians," *New England Journal of Medicine*, 327 (November 21), 1502-1506.
57. *Wall Street Journal* (1992), "Sears is Dealt a Harsh Lesson by States," October 2.
58. *Wall Street Journal* (1992), "Sears's Brennan Accepts Blame for Auto Flap," June 23, B1.
59. *Wall Street Journal* (1993), "More Managed Health-Care Systems Use Incentive Pay to Reward 'Best' Doctors," January 25, B1.
60. Weigelt, Keith and Colin Camerer (1989), "Reputation and Corporate Strategy: A Review of Recent Theory and Applications," *Strategic Management Journal*, 9, 443-454.
61. Wiener, Joshua Lyle (1985), "Are Warranties Accurate Signals of Product Reliability?" *Journal of Consumer Research*, 12 (September), 245-250.
62. Zeithaml, Valarie A. (1981), "How Consumer Evaluation Processes Differ Between Goods and Services," in J.H. Donnelly and W.R. George (eds.), *Marketing of Services*, 191-199, Chicago: American Marketing Association.
63. Zeithaml, Valarie A., Leonard L. Berry and A. Parasuraman (1988), "Communication and Control Processes in the Delivery of Service Quality," *Journal of Marketing*, 52 (April), 35-48.
64. Zeithaml, Valarie, A., A. Parasuraman and Leonard L. Berry (1985), "Problems and Strategies in Services Marketing," *Journal of Marketing*, 49 (Spring), 33-46.

# AGENCY THEORY

Susan P. Shapiro

*American Bar Foundation, Chicago, Illinois 60611; email: sshapiro@abfn.org*

■ **Abstract**   In an agency relationship, one party acts on behalf of another. It is curious that a concept that could not be more profoundly sociological does not have a niche in the sociological literature. This essay begins with the economics paradigm of agency theory, which casts a very long shadow over the social sciences, and then traces how these ideas diffuse to and are transformed (if at all) in the scholarship produced in business schools, political science, law, and sociology. I cut a swathe through the social fabric where agency relationships are especially prevalent and examine some of the institutions, roles, forms of social organization, deviance, and strategies of social control that deliver agency and respond to its vulnerabilities, and I consider their impact. Finally, I suggest how sociology might make better use of and contribute to agency theory.

## INTRODUCTION

Let me introduce myself. I am an agent. The editors of the *Annual Review of Sociology* delegated to me the task of writing an essay on agency theory. They are the principals and together we are bound in a principal-agent relationship. They have a principal-agent relationship with you (the readers) as well. They are your agents, and so am I, although not every agency theorist would agree with my loose conceptualization of your role in this, and few would be interested in you at all (although I am).

I am not sure how or why my principals selected me for this task. Perhaps they "Googled" me. I do use the words "agent," "principal," and "agency relationship" a lot. But I doubt that they used a more sophisticated search engine. If they had, they would have realized that I have never used the words "agency" and "theory" side by side (although I guess it's possible that they did and wanted someone who is not so identified with this peculiar way of understanding social reality or is not solidly in one camp or another in a rather contentious literature).

In any event, in selecting me and all the other authors in this volume, they faced a classic agency problem of asymmetric information. We know far more about ourselves—our abilities, expertise, honesty, etc.—than they do, and we sometimes make matters worse by exaggerating our talents. I know how much

**263**

of the agency literature I have bothered to read and how much of it I understand. I know whether I skip the paragraphs in the economics articles that begin, "let gamma be…" and then go on to use mathematical fonts I can't even find on my computer. I know better how good a sociologist I am and how analytical and original I am or am capable of being. I know better how many other projects I have on my plate right now and how responsible, conscientious, and diligent I am. Actually I know who would have been a better choice to write this essay. But my editors/principals don't. They never do, and therefore every assignment in this volume is tainted by adverse selection (in the insurance vernacular) or what Arrow (1985) calls "hidden information": they "will attract a disproportionate number of low-quality applicants" (Moe 1984, p. 755). The principals probably could have found someone better but just didn't know enough to identify them or didn't provide incentives compelling enough to attract them. So they got us.

Of course, that is not their only agency problem. Information asymmetries not only mean that principals don't know the true "type" (to borrow from the agency theory jargon) of the varied candidates in the pool of potential agents, but they also don't know what we are doing once they select us. They don't know what I am reading, if anything, or whether I am scouring literature reviews or plodding through the actual primary sources. They don't know whether I have been thorough or fair. They don't know if I got someone else to write this for me or if I plagiarized it. Agency theorists are mostly worried that I might be shirking—not working hard enough, if at all. Many theorists also assume that I am "opportunistic" [pursuing self-interest with "guile" (Williamson 1975)] and will take advantage of the "perquisites" of this appointment for my own benefit. But sadly, my agency-savvy principals didn't give me any perquisites. (I have tried to use my inside information to trade on *Annual Review* futures, but I can't find this product on any of the commodities exchanges.) My principals, then, are also threatened by the version of informational asymmetry known in insurance as moral hazard, or what Arrow (1985) labels "hidden action."

The one thing they can be sure of is that our goals are incompatible. My principals want the "highest-quality scientific literature reviews in the world" that "defin[e] the current state of scientific knowledge," and they want them on time and in the correct format (Annual Reviews 2003, pp. 2, 18). I want the glory with none of the work and desperately need the deadline to be extended. And I will exploit all the information asymmetries I can contrive to insure that I maximize my own interests at their expense.

So what do the poor principals do? Agency theory dictates that my principals will try to bridge the informational asymmetries by installing information systems and monitoring me. My manuscript will be peer reviewed, for example. And because my reputation is tied up in the quality of my work, they can count on some self-regulation on my part. They also offer me incentives in an effort to align my interests with theirs. They tell me that the earlier my manuscript arrives, the closer it will be placed toward the front of the volume. [So the position of this

chapter tells you something about my character, that is, if my principals are of the trustworthy type—something the sociologists (Perrow 1986), but apparently not the economists, are worried about.]

As part of this incentive alignment, my principals compensate me, not for my agreement to do this work for them or for the amount of time I spent on the project—consistent with a "behavior-oriented contract"—but for what I actually deliver, an "outcome-oriented contract." They tell me that if the manuscript arrives late, they will not guarantee that they will publish it at all, ever (and you know how difficult it will be to recycle this sort of review essay into another journal). That, of course, shifts the risk to me, because events outside of my control (like the fact that a lightening strike or virus fried my hard drive) or other environmental uncertainties may affect my ability to deliver on our agreement. Agency theory reminds us that, although principals are risk neutral (they have diversified and have plenty of other manuscripts to use), agents are risk averse, because they have placed all their eggs in this one basket. That is another reason our interests conflict, by the way; shrouded behind my information asymmetries, I will do perverse things contrary to my principal's welfare to protect myself from risk. All these efforts undertaken by my principals, coupled with the fact that I still didn't give them exactly what they wanted, constitute agency costs. The trick, in structuring a principal-agent relationship, is to minimize them.

This introduction more or less represents a cartoon version of the classic economics account of agency theory. I begin here because, as in many things, the economics formulation of agency theory is the dominant one and casts a very long shadow over the other social sciences. Because it gets all the attention and there are already excellent reviews of this literature (e.g., Moe 1984; Eisenhardt 1989; Mitnick 1992, 1998), this essay briefly traces some of the alternative disciplinary approaches—especially in law, management, and political science. Then I turn to sociology, where the literature on agency theory is especially sparse, and ask how it could be that a relationship—acting on behalf of another—that could not be more profoundly sociological does not seem to have a niche. Finally, I suggest what that niche might look like.

## ECONOMICS AND BEYOND

The main thing missing from my cartoon version of the economics of agency (unfortunately, from this agent's perspective) is any money changing hands. Consequently, a few of the traditional options for aligning my incentives with my principals (commissions, bonuses, piece rates, equity ownership, stock options, profit sharing, sharecropping, deductibles, etc.) are missing, as are some of the governance mechanisms or devices principals contrive to monitor their agents (e.g., boards of directors, auditors, supervisors, structural arrangements, and so forth). Also missing are a few of the things I might have done to reassure my principals or keep their monitors at bay: I could have bonded myself or perhaps posted

a hostage who or which wouldn't be released until I turned over the manuscript. All these devices also figure into the accounting of agency costs.

Nonetheless, my case study actually accords better with classic agency theory in economics than the scenarios economists usually model. Ours is a dyadic relationship between individuals; economists study firms and typically focus on the relationship between owners and managers or employers and employees. The assumption of methodological individualism makes this transformation seamless. In the classic articulation of agency theory in economics, Jensen & Meckling (1976) assert that "most organizations are simply legal fictions which serve as a nexus for a set of contracting relationships among individuals" (p. 310). In this paradigm, agency relationships are contracts, and the incentives, monitoring devices, bonding, and other forms of social control undertaken to minimize agency costs constitute the elements of the contract.

Economists make problematic the nature of these contracts. Those with a mathematical bent (in what is known as principal-agent theory) model the "structure of the preferences of the parties," "the nature of uncertainty," and "the informational structure" on contracting practices. A more descriptive and empirical trajectory (known as positive agency theory) examines "the effects of additional aspects of the contracting environment and the technology of monitoring and bonding on the form of the contracts and the organizations that survive" (Jensen 1983, p. 334; see also Eisenhardt 1989).

The assumption that complex organizational structures and networks can be reduced to dyads of individuals is one of many assumptions—regarding efficiency and equilibrium, that individuals are rational and self-interested utility maximizers prone to opportunism, etc.—that are off-putting to other social sciences. To be tractable, however, mathematical modeling requires such simplistic assumptions, as even a very flattering review of that literature concedes:

> [S]uch a framework sometimes encourages highly complex mathematical treatment of trivial problems; form tends to triumph over substance, and analytical concerns tend to take on lives of their own that have little to do with the explanation of empirical phenomena. . .. [M]uch of the current literature focuses on matters of little substantive interest (Moe 1984, p. 757).

One of the economists most identified with agency theory admits that "authors are led to assume the problem away or to define sterile 'toy' problems that are mathematically tractable" (Jensen 1983, p. 333).

Much of the scholarship on agency outside of economics begins by relaxing or jettisoning the unrealistic assumptions of the economics paradigm and transforming the rigid dichotomies into more complex variables. The first assumption to go, of course, is that of a dyadic relationship between individuals. As Kiser (1999) observes, classic agency theory "is an organizational theory without organizations" (p. 150). Scholarship across many disciplines brings organizations of all sorts back in and looks far beyond the economists' favorite poster children of shareholder/manager and employer/employee as they investigate when and how

agency relationships are established and regulated. Looking beyond the abstract, cloistered dyad also reveals that actors are not just principals or agents, but often both at the same time—even in the same transaction or hierarchical structure. I may be an agent to the editors of the *Annual Review*, but I am also the would-be principal to the scores of research assistants who I wish existed to assist me on this project. The CEO may be an agent of the stockholders and the board, but he or she is simultaneously the principal in a long chain of principal-agent relationships both inside and outside the corporation. What occurs at some node in that network of agents acting on behalf of the CEO figures significantly in the agency contract between the CEO and the shareholders. Just ask Kenneth Lay at Enron.

Moreover, the assumption of a solitary principal and agent is invariably extended to include multiple principals and agents. This is not just a matter of verisimilitude. Theories become much more complex (and interesting) when they allow for the possibility that collections or teams of principals (or agents) disagree or compete over interests and goals—a feature of agency relationships Adams (1996) dubs the "Hydra factor." How do agents understand and reconcile the duties delegated to them when they are receiving mixed messages and conflicting instructions—and incentives—from multiple principals? How do they do so when the contract is exceptionally vague by design, to paper over the irreconcilable differences among principals with conflicting interests—say, controversial legislation that requires implementation? When do these cleavages among and collective action problems faced by principals give agents opportunities to play one principal off against another?

Multiple agents who have been delegated to undertake a task collectively add other wrinkles to the economists' models. Agents, too, have competing interests; indeed the interests of some agents may be more congruent with those of their principals than with the other agents. Some agents are more risk averse than others; incentives work differently on different agents. Some agents may be free riders. And the existence of multiple principals and multiple agents sometimes increases the informational asymmetries and the difficulties of monitoring. These asymmetries are among the reasons organizational crimes can flourish undiscovered for long periods of time buried in complex structures of action. At other times, multiple parties help to right the imbalance of information, such as when competitive agents leak information to principals in an effort to get an upper hand over other agents (Waterman & Meier 1998).

The assumption that principals are in the driver's seat—specifying preferences, creating incentives, and making contracts that agents must follow—is also problematic (Heimer & Staffen 1998, Sharma 1997). When principals seek out agents for their expert knowledge, when principals are one-shotters and agents repeat players, when principals are unexpectedly foisted into a new role with no time or life experience to formulate preferences, let alone a contract or monitoring strategy [e.g., the new parents of a critically ill newborn (Heimer & Staffen 1998)], the asymmetry of power shifts from the principal to the agent.

Other scholars remove the economists' blinders that cause them to focus only on the self-interest and opportunism of agents and the difficulties of regulating them. Perrow (1986), for example, accuses the economics paradigm of being incapable of keeping its eye on both sides of the principal-agent relationship and of recognizing that agency problems on the agent side of the relationship are often mirrored on the principal side. He observes that the theory is indifferent to principal type that may lead to adverse selection by agents who may be unwittingly drawn to principals who shirk, cheat, and opportunistically seize perquisites for their own use; who deceive (e.g., about hazardous working conditions, opportunities for advancement, etc.); and who exploit their agents. Blind to the asymmetries of power that course through these relationships, classic agency theory, Perrow argues, is profoundly conservative, even dangerous.

Perrow (1986) also rejects the assumption that parties are invariably work averse, self-interested utility maximizers. He observes that in some settings or organizational structures, human beings are other-regarding, even altruistic, and he faults classical agency theory for its inattention to the cooperative aspects of social life. This critique is continued in what has become known in the management literature as stewardship theory, which views agents as good stewards and team players and replaces assumptions of opportunism and conflict of interest with those of cooperation and coordination (Donaldson 1990).

As other disciplines wander away from the market as the site of theoretical and empirical work on agency, the irrelevance or variability of the classic assumptions and solutions to the agency problem becomes even more apparent (Banfield 1975). Work in political science particularly confronts the limitations of a theory of markets. As Moe (1984) observes,

> the more general principal-agent models of hierarchical control have shown that, under a range of conditions, the principal's optimal incentive structure for the agent is one in which the latter receives some share of the residual in payment for his efforts, thus giving him a direct stake in the outcome. . . . For public bureaucracy, however, there is no residual in the ordinary sense of the term (p. 763).

There is no profit that can be distributed to members of public agencies for exemplary behavior. Scholarship on agency relationships, such as between the legislative or executive branch and administrative agencies, may continue to employ economic metaphors: Politicians need to maximize their votes; bureaucrats need to maximize their budgets. But the metaphor fails to capture the range of incentives at play in the political arena, many of which revolve around policy rather than profit (Waterman & Meier 1998). Indeed, the salience of policy commitments undermines our expectation of goal conflict between principals and agents, who may sometimes share policy goals (or, more accurately, some among the collections of multiple principles and agents might do so). The extent, sources, and strategies of compensating for information asymmetries also vary considerably as one moves away from market settings (Waterman & Meier 1998, Worsham et al. 1997, Sharma 1997, Banfield 1975).

Finally, scholars from varied disciplines outside of economics also abandon the assumption of an acontextual, ahistorical, and static relationship between principals and agents (Mitnick 1992). Agency relationships are enacted in a broader social context and buffeted by outside forces—other agency relationships, competitors, interest groups, regulators, legal rules, and the like—that sometimes right informational imbalances, offer or constrain incentives, exacerbate the risk of adverse selection or moral hazard, provide cover or opportunity for opportunism, and so forth. Relationships endure over time, affording principals and agents occasions to gather data about one another. Principals learn better which incentives are likely to work. Agents learn more about the preferences of the principals they serve. They develop reputations. Relationships become embedded as parties develop histories and personal relationships and become entangled in social networks (Granovetter 1985). Over time, agents acquire constituencies other than their principals that buffer them from the contracting, recontracting, and sanctioning of their principals. And as agents (government bureaucrats, corporate managers) outlast their principals (legislators, CEOs), the balance of power between principal and agent may shift.

## Management

The agency theory paradigm, first formulated in the academic economics literature in the early 1970s (Ross 1973, Jensen & Meckling 1976) had diffused into the business schools, the management literature, specialized academic and applied practitioner journals, the business press, even corporate proxy statements by the early 1990s, representing a new zeitgeist and becoming the dominant institutional logic of corporate governance (Zajac & Westphal 2004). Corporations announced the adoption of new policies, explicitly invoking agency theory buzzwords about aligning incentives, discouraging self-interested behavior by managers, and reducing agency costs. Indeed, some adopted new policies that embraced an agency rationale without bothering to implement them, simply jumping on the bandwagon of a socially constructed institutional logic that bestowed increased market value on symbolic declarations alone (Zajac & Westphal 2004).

Despite the fascinating case study in social movements (Davis & Thompson 1994), the diffusion of innovations, and the sociology of knowledge that these developments offer, they also had a significant impact on the intellectual agenda of the academy, spawning a massive empirical literature in management and organizational behavior. Agency theory has become a cottage industry that explores every permutation and combination of agency experience in the corporate form. Because the work is largely empirical, it by necessity relaxes some of the assumptions of classic agency theory in economics; it turns dichotomies into continuous variables, breathes life into abstract categories, and situates inquiry in at least some limited context. Still, it is closely wedded to the questions raised in economics and the settings invoked by economic models.

The most popular stream of literature focuses on incentive alignment, particularly compensation policies. Empirical studies consider the types and correlates

of and trade-offs between behavior-oriented (salary) and outcome-oriented (piece rates, commissions, bonuses, equity ownership and other devices that link compensation to shareholder wealth) compensation (Eisenhardt 1989). A second stream examines corporate governance and control, such as

- the monitoring role of the board of directors and trade-offs between recruiting inside or outside directors or between separating the roles of board chair and CEO versus filling them with one individual;
- monitoring strategies within the firm [e.g., trade-offs between horizontal (peer-to-peer) and vertical (agent-to-principal) control];
- bonding mechanisms; and
- the agency implications of different forms of capitalization (e.g., paying out dividends and thereby limiting discretionary funds available to managers while also activating the monitoring role of the financial markets when managers must solicit additional funding).

The literature also includes studies of the process and costs of searching for agents, especially in light of the tensions posed by adverse selection.

Another major body of scholarship considers the agency problems, agency costs, efficacy, and trade-offs of different control mechanisms as they intersect and vary by

- length of principal-agent relationship;
- organizational structure and form (e.g., headquarters and subsidiary, outsourcing);
- characteristics of industries, organizations, and employees (e.g., technologies, product demand, diversification, venture capitalist-entrepreneur relationships, family firms, cultural distance between sites, employee education, skill levels, amount of specialized knowledge, autonomy, etc.);
- "programmability" of the task, or how well the required behaviors can be precisely defined (Eisenhardt 1989); and
- organizational environments (e.g., turbulence).

Also coursing through this literature is a debate, sketched earlier, between those who adopt the skeptical, even paranoid, assumptions of agency theory and the costly control mechanisms it propounds and those who have a more hopeful view of human capacities for other-regarding behavior and altruism and argue that agency costs can be mitigated by organizational structures that foster reciprocity, cooperation, embeddedness, and trust (Donaldson 1990, Wright & Mukherji 1999).

## Political Science

In exploring the delegation of power and authority in political and government institutions and international organizations, political scientists take agency theory outside of the economic marketplace and the constricting web of assumptions

that shroud the economic theory of agency. The political system can, of course, be understood as a complex network of principal-agent relationships composed of citizens, nation states, elected officials, lawmakers, members of the executive branch, administrative agencies, courts, international organizations, ambassadors, bureaucrats, soldiers, police officers, supervisory officials, civil servants, patronage appointees, and even those who monitor other agency relationships inside political institutions and in the market. These actors concurrently play principal and agent roles within and across political organizations.

A general theory of agency emerged in political science (Mitnick 1973) at the same time that it did in economics (Ross 1973), apparently independently. As we have seen, the latter took off spectacularly, becoming quickly institutionalized in an academic literature, specialty journals, and corporate ideologies and practices. The former languished (Moe 1984), developing belatedly as rational choice theory made inroads into political science. As a result, agency theory in political science borrows heavily from the economics paradigm rather than the more sociological conception offered by Mitnick (1973) or even classic works, such as Weber on bureaucracy (Kiser 1999).

The vague outlines of the agency paradigm in political science are the same as those in the classic version: Principals delegate to agents the authority to carry out their political preferences. However, the goals of principals and agents may conflict and, because of asymmetries of information, principals cannot be sure that agents are carrying out their will. Political principals also face problems of adverse selection, moral hazard, and agent opportunism. So principals contrive incentives to align agent interests with their own and undertake monitoring of agent behavior, activities that create agency costs.

The details are quite different, however, for many of the reasons considered earlier. Political scientists assume multiple agents and principals; heterogeneous preferences or goal conflict and competition among principals and among agents as well as between them; problems of collective action; a more complicated palate of interests and therefore different incentives mobilized to control them; varying sources of and mechanisms to mitigate informational asymmetries; an active role for third parties (interest groups, regulated parties, etc.); and a dynamic playing field on which relationships unfold and are transformed.

Political scientists also consider a more diverse set of scenarios for delegating power beyond those inherited from the economics paradigm. Principals may delegate to another to enhance the credibility of their commitments, for self-binding (to ensure their long-term resolve in the face of immediate temptations), or to avoid blame for unpopular policies. These scenarios call for a very different agency contract. Instead of providing incentives and sanctions to align the interests of agents with their own, principals seeking credibility from their agents select agents operating at arm's length, with very different policy preferences, and confer considerable discretion and autonomy to them. These agency contracts grant independence while still seeking to insure accountability (Majone 2001).

Early literature in political science on the iron law of oligarchy, the iron triangle (between Congress, regulatory agencies, and regulated interests), regulatory

capture, and bureaucratic drift all give voice to some of the intrinsic difficulties of principal control in political institutions. More recent work employing an agency theory perspective ranges from appellate review of lower court decisions to political corruption and presidential decisions to use force. The largest literatures examine state policy implementation, the relationship between elective institutions and administrative agencies (especially legislators and bureaucrats), and government regulation. Principal-agent perspectives are also commonplace in examinations of international organizations (e.g., central banks, international courts, the European Union) in the literature on comparative politics and international relations.

Political scientists devote far more attention than economists to the details of how principals control agents. There is some work on the selection and recruitment of agents, the role of patronage, political appointments, and the impact of civil service requirements on adverse selection and more on how principals specify their preferences. A body of work considers statutory control (i.e., detailed legislation) and how lawmakers craft legislation to restrict the discretion of those charged with its implementation, specifying administrative structures and procedures to constrain the decision-making process (McCubbins et al. 1989). There are literatures on political oversight and monitoring, including ways in which principals opt for reactive over proactive oversight, relying on third-party monitoring by affected interest groups or the targets of their legislation to detect and report on noncompliance (so-called fire alarms or decibel meters).

There is more attention in political science than in economics to the role of sanctions—budget cuts, vetoing rules or agency actions, reversing court decisions, firing officials or voting them out of office, requiring agency reauthorization or threatening recontracting, etc.—perhaps because, as noted earlier, it is far less easy to align incentives without the financial inducements that flow through economic organizations. The literature also considers the matter of agency costs; when they are too high, principals may decide not to squander resources on them (Mitnick 1998, Banfield 1975). Because politicians may not directly feel the consequences of self-interested, opportunistic agents shirking or undermining their interests (what political scientists call slack, slippage, or bureaucratic drift), the costs of which are generally passed along to the public, monitoring activities may be more lax in political arenas (Waterman & Meier 1998).

## Law

Long before there was a theory of agency, there was a law of agency. Indeed, it was not until the twenty-first century that the *Restatement of the Law, Agency* (American Law Institute 2001) replaced "master/servant" with "employer/employee." The law of agency

> encompasses the legal consequences of consensual relationships in which one person (the 'principal') manifests assent that another person (the 'agent') shall, subject to the principal's right of control, have power to affect the principal's legal relations through the agent's acts and on the principal's behalf (American Law Institute 2001, p. 1).

In other words, the central focus of the law of agency is on "the legal consequences of choosing to act through another person in lieu of oneself" (DeMott 1998, p. 1039). Agency doctrine defines the legal obligations that principals have with third parties for actions that agents took on their behalf. The principal, for example, may be "bound to contracts and transactions made by the agent and may be vicariously liable for some instances of the agent's misconduct" (DeMott 1998, p. 1038). Because principals will be held responsible for the actions of their agents, the law also attends to the sources of agent authority, clearly demarcating what constitutes an agency relationship, the rights of principals to control their agents, and the fiduciary duty and other obligations that agents owe their principals (Clark 1985).

Agency theory borrows jargon from agency law, but adopts neither its definition nor its central focus. The legal definition of agency is much more narrow even than that employed in the economics paradigm of agency theory, let alone that found in the other social sciences.

> [A]gency does not encompass situations in which the 'agent' is not subject to a right of control in the person who benefits from or whose interests are affected by the agent's acts, who lacks the power to terminate the 'agent's' representation, or who has not consented to the representation (American Law Institute 2001, p. 2).

> Generally, the alleged agent and principal have met each other face to face, or have talked on the telephone, or have otherwise communicated with each other in a specific, individualized way. Courts trying to determine the scope of their relationship often scrutinize the actual course of dealings between the particular parties and try to determine what their actual understanding of their particular relationship was (Clark 1985, p. 58).

The relationship between a corporation's shareholders and its directors, for example, does not fall within the legal definition of agency, notwithstanding the centrality of this relationship in economic agency theory. Principal control is critical in the law of agency because of its focus on third parties and the concern that when third parties make agreements with agents or are hurt by agents, their principals will be bound or held responsible. But it is the control itself that the social sciences make problematic. Therefore, it cannot be defined away by looking only at the point along a continuum where control is absolute. Moreover, central questions in the social sciences about the nature of the contract between principal and agent, the mechanisms by which the former control the latter, and strategies to contain agency costs are rather peripheral in the law of agency.

Still, when the two paradigms do intersect, the law of agency provides rich grist for the social scientists' mill—for example, when legal scholars look to the mechanisms by which principals select their agents; the private norms, instructions, and messages the principals convey; the nature of the incentives they offer; and the care they take to monitor the behavior of agents to determine whether corporations should be held vicariously liable for the criminal conduct of their employees (DeMott 1997). The law offers normative understandings of agency relationships

and lots of data (if tainted by selection bias), especially when they fail. But it offers little else.

## A SOCIOLOGICAL PERSPECTIVE

> Although economists may speak of 'the agency problem,' agency is in fact a solution, a neat kind of social plumbing. The problem is the ancient and ineluctable one of how to attain and maintain control in order to carry out definite, yet varying purposes (White 1985, p. 188).

In his comparative analysis of agency theory applications to state policy implementation in economics, political science, and sociology, Kiser (1999) observes that, compared to the other two disciplines, "the use of agency theory in sociology is in its infancy" and comes from a rather different "intellectual genealogy" (p. 162), largely the work of Weber (1924/1968). [See Kiser (1999) for an illuminating analysis that traces the linkages between abstract components of classic agency theory and Weber's work on the relationship between rulers and their administrative staff.]

Empirical work in sociology that explicitly adopts an agency theory perspective (aside from that described earlier in the organizational behavior and management literatures) can be found in the most unexpected of places—in qualitative comparative historical sociology. In imaginative and richly textured case studies of such things as European colonialism in seventeenth and eighteenth century Asia, Chinese state bureaucratization that occurred two millennia before any of the European states, early modern tax farming, and types of corruption in premodern Asian tax administration, we learn about the tensions between principals and agents, conflicting interests, opportunism, informational asymmetry, agent selection, monitoring, sanctions, incentives, and agency costs (Adams 1996, Kiser 1999, Kiser & Cai 2003). This work links social structure to types of agency relations, and it demonstrates how different combinations of recruitment, monitoring, and sanctioning practices yield different administrative systems. This literature is certainly a far cry from the abstract mathematical models of principal-agent theory in economics.

It is puzzling that agency theory is not invoked elsewhere across the sociological landscape in places one would think would be more hospitable. Perhaps, like me, few sociologists feel comfortable putting the words "agency" and "theory" side by side and find the classic paradigm, its assumptions, and the research questions it inspires off-putting and simplistic. But that has never been our only choice. As long as there has been an economic theory of agency there has been a more sociological alternative. In a series of papers spanning at least 25 years, political scientist Barry Mitnick broke the monopoly on agency theory enjoyed by the economics paradigm and offered an alternative to the assorted baggage that comes with it. Agency, he argued (Mitnick 1998, p. 12) is simply "a general social theory of relationships

of 'acting for' or control in complex systems." Agency relationships have two faces, Mitnick observed: "the activities and problems of identifying and providing services of 'acting for' (the agent side), and the activities and problems of guiding and correcting agent actions (the principal side)." Of course, both faces of agency entail costs and at some point it does not pay for principals or agents to perfect their behaviors. So "perfect agency" is rare, and deviant behavior is likely to "persist and be tolerated." Agency theory, then, "becomes a study in the production, the persistence, and the amelioration of failures in service and control" (Mitnick 1998, p. 12), a kind of Murphy's law (Mitnick 1992, p. 76). Mitnick's work repeatedly shows the links between agency theory and sociological literatures from exchange theory to norms, networks, authority, organizations, social control, regulation, trust, social cognition, and so on. Yet it, too, is rarely cited in sociological literature.

The problem may be that "acting for" relationships are too general, embracing too much of what is enacted on our turf. Perhaps sociologists have been studying agency all along and just didn't know it. In the remainder of this essay, I focus on several sites across the social landscape where making agency relationships problematic seems likely to provide the most theoretical purchase.

Agency or "acting for" relationships arise from a number of sources, including

1. the division of labor; we simply do not have time to do everything ourselves (even hunting and gathering), and complex tasks often require more than one actor [Mitnick (1984) calls this practical or structural agency];

2. the acquisition of expertise or access to specialized knowledge [Mitnick (1984) labels this contentful agency];

3. the bridging of physical, social (e.g., brokering or intermediation), or temporal distance [Adams's (1996) study of colonialization provides an example of the challenges of the former; for the latter, see Majone's (2001) discussion of time-inconsistency]; and

4. the impulse to collectivize in order to enjoy economies of scope and scale or protection from risk [Mitnick (1984) calls this systemic or collective agency]; many of these relationships (pensions, insurance, investments, etc.) are what I have called futures transactions that "demand that commitment be conferred far in advance of payoff without any necessary confirmation during the interim that the return on investment will ever be honored" (Shapiro 1987, p. 628).

These varied occasions for agency—especially the last three, in which a formidable physical, social, temporal, or experiential barrier separates principal and agent—pose different agency problems. Several exacerbate problems of asymmetric information; others contribute to adverse selection; some create collective action problems among multiple principals; others provide easy cover for moral hazard and opportunism.

## Professions

The sociology of the professions provides a window on agency as expertise, problems of asymmetric information, and one kind of model for delivering agency services. The assumptions of the agency paradigm are stretched where principals seek out agents for their specialized knowledge. Sharma (1997) observes that run-of-the-mill information asymmetry (not knowing what the agent does) is exacerbated in encounters with professionals by knowledge asymmetry as well (not knowing how the agent does a job). Adverse selection is a special problem because principals are unable to evaluate the skills of prospective agents. Principals also have a difficult time specifying an agency contract because they may not know what expert services are required or how much of them, what procedures ought to be followed, or what criteria are appropriate to limit agent discretion. They also have difficulty evaluating the quality of service because "indeterminacy [is] intrinsic in highly specialized tasks" (Sharma 1997, p. 771). Some patients get better despite their physicians; the clients of superb lawyers sometimes lose; and bright, curious, conscientious students may become great sociologists despite incompetent or opportunistic professors.

Professions provide the solution to these agency problems. They boast careful and competitive selection procedures. They offer training and credentialing, licensing, recertification, and mandatory continuing education to solve the principals' problem of adverse selection. They may even establish protocols or specify best practices to limit agent discretion. They create ethics codes to curb the self-interest and opportunism of practitioners. Because principals are unable to determine when they have received exceptional or substandard service, professions self-regulate in varied settings (among peers, within service organizations, within professional associations, and by disciplinary bodies). And professions often offer or promote malpractice insurance to protect principals from the errors or misdeeds of honest and incompetent agents alike. Insurers often provide incentives, stipulate mandatory procedures, and provide loss prevention services to their insureds—adding yet another level of regulation (Heimer 1985, Davis 1996). Professions, then, are social devices to limit agency costs.

Of course, there is a critical literature that provides a rather different frame on the agendas of professions as mechanisms to secure monopoly (e.g., Larson 1977, among many others). But this frame is by no means incompatible with a principal-agent perspective. Indeed offering a credible mechanism to minimize agency costs represents a brilliant marketing strategy and a way to stave off the encroachment of other would-be agents who seek to offer the same services to principals.

## Embeddedness

Literatures on embeddedness and trust (Granovetter 1985, Shapiro 1987, Cook 2001) depict a rather different strategy for coping with the agency problem by targeting agent selection, monitoring, and sanctioning. Embedding agency relationships in an ongoing structure of personal relationships solves the problem of

adverse selection in the recruitment of agents. Principals frequently know their agent's type because of personal familiarity with potential agents or through members of trusted social networks in which both principal and agent are embedded; agents have track records and reputations. Although neither self-interest nor goal conflict is extinguished by recruiting agents from personal networks, their effects are likely mitigated somewhat. Agents and principals are more likely to share similar interests and values than those found among groups of strangers, and agents are more likely to be other-regarding (altruistic, even) or honest when entrusted with responsibilities for friends, family, neighbors, fellow church or association members, and the like. Monitoring of agent behavior is also usually easier in proximate and continuing relationships in which agents are routinely overseen or surveilled by principals or their associates. And social networks afford a rich array of sanctions for the errant agent (from shaming, ostracizing, or loss of reputation, to more restitutive sanctions).

Despite the celebration of trust as a source of social capital in the literature, embeddedness also has a dark side. Family firms, for example, face unique agency costs. They struggle with adverse selection because nepotism can lead to the selection of less-capable or expert agents. Moreover, because family members are often compensated generously regardless of merit, and their job tenures are relatively secure, principals lack important incentives to constrain agent behavior. Hence, the risk of shirking and free riding by family agents. Because embeddedness is often an excuse to relax vigilant recruitment and monitoring, it provides cover, not only for wayward offspring or relatives, but also for confidence swindlers to feign social intimacy and thereby enjoy unfettered opportunism (Shapiro 1990).

## Fiduciaries

In the law of agency, all agents are fiduciaries, but all fiduciaries are not agents (that is because, as you recall, in law agents must be able to control their principals). But these other non-agent fiduciaries are much more interesting—the individuals and organizations acting on behalf of those for whom the asymmetries of information, expertise, access, or power are so great that they cannot pretend to control their agents. We are more interested in the professor who has his pension tied up in TIAA-CREF than the CEO of TIAA-CREF who has delegated some responsibility to an investment analyst working at the company. We are more interested in Terry Schiavo, the comatose Florida woman whose guardian is trying to end life support, than in Jeb Bush, the Florida governor who is maneuvering to continue her persistent vegetative state. Or, more accurately, I propose that sociologists take an interest in the fiduciaries acting on behalf of the former. Organizational and political sociologists have already taken an interest in the agents for the latter.

When agency relationships are at their most asymmetric, the basic logic of classic agency theory breaks down. Preferences are not specified (or at least not heard or satisfied), contracts not formulated, incentives not fashioned, monitoring not mobilized, sanctions not levied—at least not by the principals themselves;

and those who believe that agents are opportunistic might profitably look here for evidence of abuse. Of course, these fiduciaries face a problem as well: Why would anyone ever trust them when their conduct is so unrestrained? Would-be fiduciaries therefore undertake activities to shore up their trustworthiness in an effort to market their wares. The systematic study of the social construction, social organization, and social control of the fiduciary role or impersonal trust is well overdue [Shapiro (1987); see also Majone (2001) for a discussion of trustee or fiduciary relations as an alternative to agency in political science].

## Goal Conflict

The classic agency paradigm, with its eye on the principal, perceives goal conflict as the departure of agents from the interests of the principal. Hence, the solution to this agency problem is to come up with incentives that will align the interests of agents with those of the principal. Keep the agent from shirking by paying her a piece rate, perhaps. The agency problem looks quite different from the perspective of the agent, though. Conflicts between the interests of the agents and those of the principal are the least of the agent's problems. The real problem is that the agent is most likely serving many masters, many of them with conflicting interests. Even if the agent is able to silence his or her own interests, there is the matter of how to maneuver through the tangled loyalties he or she owes to many different principals and how to negotiate through their competing interests and sometimes irreconcilable differences. How do you honor the preferences of one when doing so means that you are undermining the interests of another? Can you represent a client suing an insurance company if another lawyer in your firm represents insurance companies? Do you take your patient off antipsychotic drugs because your clinical trial requires subjects begin with a drug washout (possibly followed by a placebo)? Do you audit a company that pays your firm millions of dollars annually for management consulting services? Do you take the kidney of one of your offspring to save another offspring, or perhaps conceive one to use its stem cells or bone marrow for another? Do you read the dissertation or peer review the article? How do agents choose among often incommensurable interests that do not share a common metric along which competing demands can be ranked, costs and benefits weighed, trade-offs evaluated, or rational choices modeled (Espeland & Stevens 1998)?

Only the rare agent has the luxury of aligning her interests with a single principal. Conflict of interest is hardly about shirking or opportunism with guile; it is about wrenching choices among the legitimate interests of multiple principals by agents who cannot extricate themselves from acting for so many. In an economy driven by mergers, diversification, cross-ownership, synergy, interdisciplinary practices offering one-stop shopping, and dizzying job mobility, agents are increasingly buffeted by the conflicting interests of the principals they serve. Classic agency theory misunderstands not only the source of goal conflict but also the social conditions that inflame it. Examining how the social organization of agency relationships

gives rise to conflicting interests and how agents (institutional as well as individual) in diverse settings and roles respond is a subject ripe for sociological inquiry (e.g., Shapiro 2003).

## Opportunism

Of one thing classic agency theory is sure: There will be agency problems. But it is remarkably vague about the nature of the problems, short of shirking and exploiting perquisites. The term guile does not quite spell out what agents are up to when they act opportunistically either. Sociologists have been studying these agency problems at least since Edwin Sutherland (1940) coined the term white-collar crime in his presidential address to the American Sociological Society. After many years of spirited disagreement, sociologists now agree to disagree about the appropriate definition of white-collar crime. But, aside from those who continue to insist that these are merely the crimes of high-status individuals, many would probably agree that misdeeds committed by individual or organizational agents come fairly close to what they consider to be white-collar crimes. I go further, asserting that we focus on the fiduciary duties of those in positions of trust, and I define white-collar crime as "the violation and manipulation of the norms of trust—of disclosure, disinterestedness, and role competence" (Shapiro 1990, p. 250). But I am not sure that I have convinced other sociologists. Nonetheless, few would contest the characterization of lying (misrepresentation and deception) and stealing (misappropriation, self-dealing, and corruption) by those in positions of trust (i.e., agents) as core elements of what they mean by white-collar crime. Nor would many argue that understanding how the structural properties of agency relationships facilitate misconduct and confound systems of social control is not central to agency theory models regarding policing and sanctioning of agent opportunism.

Although traditional agency theorists write frequently about corruption and probably mean misappropriation or self-dealing when they refer to the exploitation of perquisites, I doubt they would be altogether comfortable with this approach. A whistleblower, for example, would be violating the agency contract as would an employee who silently refused to be complicit in organizational misconduct ordered by his or her principals. Neither of these agency-theory malefactors would be problematic in a sociological conception because, unconstrained by assumptions of methodological individualism, sociologists can juggle many units of analysis and sites and chains of principal/agent relationships simultaneously. Although classic agency theorists seemed surprised when the world learned that their perfect incentives to align the interests of corporate executives and shareholders (giving the former stock options and equity ownership) might result in these executives contriving illicit schemes to inflate stock prices, sociologists, with our eyes on the bigger picture, surely were not. Nor are we convinced that these extraordinarily costly agency failures constitute a refutation of agency theory, as some suggest (Zajac & Westphal 2004); rather, we argue that one needs a more nuanced understanding of principals, agents, and organizations when fashioning complex

incentives. (Besides, we have been trained to be mindful of the unanticipated consequences of purposive social action.)

Sociologists have and will continue to make an important contribution to understandings of white-collar and corporate crime (Shapiro 2001). Bringing the insights of agency theory to their inquiry will push the envelope a bit further and sharpen their insights.

## Monitoring

There is, of course, an abundance of work in sociology on social control, compliance, organizational governance, policing, and sanctions that will contribute to understanding the agency paradigm. There are also more specialized literatures on the cover up of organizational misconduct and the social control in and of organizations, organizational intelligence, regulation and enforcement, and the sanctioning of white-collar or corporate offenders. These literatures demonstrate that much of what we know about the control of crime in the streets does not work so well when we seek to understand crime in the suites (i.e., agency problems). I cannot possibly review them here or even supply the dozens of citations to the most groundbreaking work in this area.

However, two observations are relevant here. First, because information and knowledge asymmetries ("know what" and "know how") are characteristic of many agency relationships, and because agency relationships are exceptionally opaque [owing to institutions of privacy (Stinchcombe 1963)] and relatively inaccessible to surveillance, self-regulation (drawing on inside information and expertise) plays an important monitoring role. Sociologists have tended to be skeptical of self-regulation—of foxes guarding chicken coops—as an institutionalized conflict of interest. Much good work has proven that stereotype simplistic (e.g., Kagan et al. 2003, Ayres & Braithwaite 1992). But, whatever the efficacy of self-regulation, it requires continued scholarly attention in the policing of agency relationships.

Second, many of the regulatory and self-regulatory arrangements devised to monitor agency relationships are themselves agency relationships. Whether they are internal or external auditors, compliance officers, internal affairs departments, government regulators, insurance companies, investment advisors, or rating agencies (e.g., Standard & Poors or Underwriters Laboratory), the monitors are acting on behalf of some set of principals. And, therefore, they too promise agency problems. They shirk, become coopted, engage in corruption, or perhaps simply monitor the wrong things. In an escalating cycle of agents overseeing agents, we must ask: Who monitors the monitors (Shapiro 1987)?

## Insurance and Risk

There is a reason that the basic language of agency theory—adverse selection and moral hazard—comes from insurance. Insurance institutions have been designing contracts and negotiating around the shoals of goal conflict, opportunism,

monitoring, and especially incentives long before the social sciences discovered agency. Insurance companies, indeed, know so much about failures of agency that they sell policies (fidelity bonds, for example, or liability policies for breach of fiduciary duty or professional malpractice) to cover such things, putting their money where their mouths are, a risk I doubt few academics would take. As Heimer (1985) demonstrated some time ago, sociologists have a great deal to learn from the social practices of insurance. They still do.

## Agency Costs

However hard principals try to minimize them, all agency relationships experience agency costs; about this all the paradigms agree. Agency costs arise from many sources: the costs of recruitment, adverse selection, specifying and discerning preferences, providing incentives, moral hazard, shirking, stealing, self-dealing, corruption, monitoring and policing, self-regulation, bonding and insurance, agents who oversee agents who oversee agents, as well as failures in these costly corrective devices. Because principals cannot observe agent behavior, they "rely on imperfect surrogate measures, which can lead the agent to displace his behavior toward the surrogates in order to appear to be behaving well" (Mitnick 1992, p. 79) (e.g., because student test scores are used to monitor teachers, some teacher/agents coach students on how to take tests rather than teaching them substance or how to think). Agency costs therefore increase because agents are concentrating their efforts on the wrong things.

Costs also increase because organizations are structured to minimize opportunism—checks and balances are created, reporting requirements implemented, redundancies introduced, employees rotated, responsibilities fragmented, layers of supervision added, revolving doors locked, and so on. Costs increase because principals, fearful of abuse, impose procedures, decision rules, protocols, or formularies to limit agent discretion, or their agents do. Ironically, principals who seek out agents because they lack the expertise to make decisions tell their agents how to make decisions on their behalf, or else they tie their hands. Although organizational sociology has demonstrated that agents sometimes bend the rules to better serve their principals, others ritualistically follow the letter rather than the spirit of the law, thereby deepening agency costs. Because we fear that agents might act on their self-interests, we require that they be disinterested; we take agents out of embedded networks where their loyalties and interests are entangled with others, but at the price of losing the social capital, reputation, goodwill, and inside information that they might have used profitably in service of their principals.

In short, because we are fearful that agents will get our preferences wrong, we construct a protective social edifice that insures that they will get them less right. As I wrote in a different context some time ago, these trade-offs between one kind of agency cost over another are akin to the choice between Type I and Type II errors in statistics. Are the constraints set so narrowly that desirable agent

behavior is deterred or so flexibly that inappropriate behavior is tolerated (Shapiro 1987)? Either way, you get an error. Mitnick (1998) reminds us that the costs are sometimes just not worth it, and perfect agency is rare indeed.

These reflections about the sources and consequences of agency costs are just that; certainly they warrant more systematic investigation. How do principals make investment decisions about agency costs? For what kinds of agency relationships are costs the highest? Aside from embedding agency service in ongoing social relationships, what strategies do principals employ to minimize agency costs? When do they simply throw up their hands and decide not to delegate at all?

## CONCLUSION

Although agency theory may not occupy a niche in sociology, agency relationships are omnipresent, under cover of other aliases—bureaucracy, organizations, professions, roles, markets, labor, government, family, trust, social exchange, and so on—"a neat kind of social plumbing," as White (1985, p. 188) observed. Drawing on agency theory in other disciplines, sociologists have been sensitized not to lose sight of the interaction between agent selection, specification of preferences, designing incentives to align the interests of principal and agent, monitoring, and sanctioning in the "acting for" relationships that unfold on their substantive terrain. But that is just the beginning. Sociology has much more to offer, as I have suggested above, both in examining the sites along the social landscape where agency is especially prominent and, having jettisoned the unrealistic assumptions and abstract models fashioned in the other social sciences, in inquiring in empirical detail about how principals and agents actually choreograph their dance. Are sociologists ready to use "agency" and "theory" side by side? I think not. But that's the good news.

**The *Annual Review of Sociology* is online at http://soc.annualreviews.org**

## LITERATURE CITED

Adams J. 1996. Principals and agents, colonialists and company men: the decay of colonial control in the Dutch East Indies. *Am. Sociol. Rev.* 61:12–28

American Law Institute. 2001. *Restatement of the Law Third, Restatement of the Law, Agency*, Tentative Draft #2 (March 14). Philadelphia: American Law Institute

Annual Reviews. 2003. *Instructions for the preparation of manuscripts*. http://www.AnnualReviews.org

Arrow KJ. 1985. The economics of agency. See Pratt & Zeckhauser 1985, pp. 37–51

Ayres I, Braithwaite J. 1992. *Responsive Regulation: Transcending the Deregulation Debate*. New York: Oxford.

Banfield EC. 1975. Corruption as a feature of governmental organization. *J. Law Econ.* 18:587–605

Clark RC. 1985. Agency costs versus fiduciary duties. See Pratt & Zeckhauser 1985, pp. 55–79

Cook KS, ed. 2001. *Trust in Society*. New York: Russell Sage Found.

Davis AE. 1996. Professional liability insurers as regulators of law practice. *Fordham Law Rev.* 65:209–32

Davis GF, Thompson TA. 1994. A social movement perspective on corporate control. *Admin. Sci. Q.* 39:141–73

DeMott DA. 1997. Organizational incentives to care about the law. *Law Contemp. Probl.* 60:39–66

DeMott DA. 1998. A revised prospectus for a third restatement of agency. *U.C. Davis Law Rev.* 31:1035–63

Donaldson L. 1990. The ethereal hand: organizational economics and management theory. *Acad. Manag. Rev.* 15:369–81

Eisenhardt KM. 1989. Agency theory: an assessment and review. *Acad. Manag. Rev.* 14:57–74

Espeland WN, Stevens ML. 1998. Commensuration as a social process. *Annu. Rev. Sociol.* 24:313–43

Granovetter M. 1985. Economic action and social structure: the problem of embeddedness. *Am. J. Sociol.* 91:481–510

Heimer CA. 1985. *Reactive Risk and Rational Action: Managing Moral Hazard in Insurance Contracts*. Berkeley: Univ. Calif. Press

Heimer CA, Staffen LR. 1998. *For the Sake of the Children: The Social Organization of Responsibility in the Hospital and the Home*. Chicago: Univ. Chicago Press

Jensen MC. 1983. Organization theory and methodology. *Account. Rev.* 58:319–39

Jensen MC, Meckling WH. 1976. Theory of the firm: managerial behavior, agency costs and ownership structure. *J. Financ. Econ.* 3:305–60

Kagan RA, Gunningham N, Thornton D. 2003. Explaining corporate environmental performance: How does regulation matter? *Law Soc. Rev.* 37:51–90

Kiser E. 1999. Comparing varieties of agency theory in economics, political science, and sociology: an illustration from state policy implementation. *Sociol. Theory* 17:146–70

Kiser E, Cai Y. 2003. War and bureaucratization in Qin China: exploring an anomalous case. *Am. Sociol. Rev.* 68:511–39

Larson MS. 1977. *The Rise of Professionalism: A Sociological Analysis*. Berkeley: Univ. Calif. Press

Majone G. 2001. Two logics of delegation: agency and fiduciary relations in EU governance. *Eur. Union Polit.* 2:103–22

McCubbins MD, Noll RG, Weingast BR. 1989. Structure and process, politics and policy: administrative arrangements and the political control of agencies. *Virginia Law Rev.* 75:431–82

Mitnick BM. 1973. *Fiduciary responsibility and public policy: the theory of agency and some consequences*. Presented at Annu. Meet. Am. Polit. Sci. Assoc., 69[th], New Orleans

Mitnick BM. 1984. *Agency problems and political institutions*. Presented at Annu. Meet. Am. Polit. Sci. Assoc., 80[th], Chicago

Mitnick BM. 1992. The theory of agency and organizational analysis. In *Ethics and Agency Theory*, ed. NE Bowie, RE Freeman, pp. 75–96. New York: Oxford Univ. Press

Mitnick BM. 1998. Agency theory. In *The Blackwell Encyclopedic Dictionary of Business Ethics*, ed. RE Freeman, PH Werhane, pp. 12–15. Malden, MA: Blackwell

Moe TM. 1984. The new economics of organization. *Am. J. Polit. Sci.* 28:739–77

Perrow C. 1986. Economic theories of organization. *Theory Soc.* 15:11–45

Pratt JW, Zeckhauser RH, eds. 1985. *Principals and Agents: The Structure of Business*. Boston: Harvard Bus. Sch. Press

Ross SA. 1973. The economic theory of agency: the principal's problem. *Am. Econ. Rev.* 63:134–39

Shapiro SP. 1987. The social control of impersonal trust. *Am. J. Sociol.* 93:623–58

Shapiro SP. 1990. Collaring the crime, not the criminal: reconsidering the concept of white-collar crime. *Am. Sociol. Rev.* 55:346–65

Shapiro SP. 2001. Crime: white-collar. In *International Encyclopedia of the Social & Behavioral Sciences*, ed. NJ Smelser, PB Baltes, 5:2941–45. Oxford: Elsevier

Shapiro SP. 2003. Bushwhacking the ethical high road: conflict of interest in the practice of law and real life. *Law Soc. Inq.* 28:87–268

Sharma A. 1997. Professional as agent: knowledge asymmetry in agency exchange. *Acad. Manag. Rev.* 22:758–98

Stinchcombe AL. 1963. Institutions of privacy in the determination of police administrative practice. *Am. J. Sociol.* 69:150–60

Sutherland EH. 1940. White-collar criminality. *Am. Sociol. Rev.* 5:1–12

Waterman RW, Meier KJ. 1998. Principal-agent models: an expansion? *J. Public Admin. Res. Theory* 8:173–202

Weber M. 1924/1968. *Economy and Society: An Interpretive Sociology*, ed. G Roth, C Wittich. New York: Bedminister

White HC. 1985. Agency as control. See Pratt & Zeckhauser 1985, pp. 187–212

Williamson OE. 1975. *Markets and Hierarchies: Analysis and Antitrust Implications*. New York/London: Free Press

Worsham J, Eisner MA, Ringquist EJ. 1997. Assessing the assumptions: a critical analysis of agency theory. *Admin. Soc.* 28:419–40

Wright P, Mukherji A. 1999. Inside the firm: socioeconomic versus agency perspectives on firm competitiveness. *J. Socio-Econ.* 28: 295–307

Zajac EJ, Westphal JD. 2004. The social construction of market value: institutionalization and learning perspectives on stock market reactions. *Am. Sociol. Rev.* 69:233–57

# A Principal-Agent Theory Approach to Public Expenditure Management Systems in Developing Countries

*by*

Luc Leruth *and* Elisabeth Paul*

*A well-functioning public expenditure management system (PEM) is considered a critical pillar of government efficiency. This article discusses PEM systems in developing countries using an analytical framework based on principal-agent theory. This simple model can be applied to various PEM systems and allows for comparisons between institutional settings. To illustrate this, the authors analyse the benefits derived from the use by the ministry of finance of* ex post *audits and* ex ante *controls, and assess their value in terms of their ability to deter cheating. The authors derive a set of possible "control regimes" which can be used by the ministry of finance.*

**1**

## 1. Introduction

A well-functioning public expenditure management (PEM) system is considered to be a critical pillar of government efficiency by most practitioners, who place it on par with a low-distortion tax system and an efficient tax administration. It is therefore unfortunate that there is so little economic research on the design of PEM systems, especially on the theoretical side.[1] On the empirical side, papers have generally focused on the efficiency of public expenditure in key sectors (health and education), and only a few attempts have been made to quantify the welfare losses associated with a weak PEM system. They all point to rather high economic costs. For example, a public spending tracking survey in Uganda concludes that only 13% of nonsalary expenditures earmarked for primary schools reached the intended beneficiaries during 1991-95. The bulk of the allocated spending was either used by public officials for purposes unrelated to education or captured for private gain (Reinikka and Svensson, 2004). In Ghana, a survey concluded that 20% of nonwage public health expenditure and 50% of nonwage education expenditure reached the frontline facilities (Ye and Canagarajah, 2002).

The importance of a good PEM system has come to the forefront of the debate in the context of the debt initiative for Heavily Indebted Poor Countries (HIPCs), which provides substantial debt relief from the international community while requiring eligible countries to pursue good economic policies and to make their budget more "pro-poor" using the HIPC relief for spending on priority areas of a country's poverty reduction strategy. The difficulty in tracking public expenditure has become clear during the systematic assessment of the capacity of some 25 HIPCs by international financial institutions.[2] Without getting into the details of the methodology used to assess PEM systems, it is worth mentioning that it was based on 15 benchmarks (extended to 16 in 2004) relating to the three main components of budget management.[3] The studies indicate that, while progress has been made since the initiative was launched, a majority of HIPCs still require substantial upgrading of their PEM systems to be capable of reliably tracking public spending. In particular, internal control and the production of final audited accounts are the areas in most need of strengthening.

The problem is that the list of recommended key reforms for "getting the basics right" (Schick, 1997) is quite large and, although internal and external controls are identified as priorities, the list of priorities also covers most other

areas (see Diamond, 2006). As HIPCs are, by definition, severely constrained in terms of both financial and human resources, it therefore appears critical to address the key areas of weaknesses in the most effective way possible.

In this article, we use the well researched principal-agent theoretical framework to clarify the issues arising in PEM systems and help prioritise the PEM reform agenda. It has already been argued that a chain of principal-agent relationships characterises PEM systems, which in turn raises the potential for agency problems (see, for example, Tanzi, 2000, p. 445).

We interpret corruption and bad governance as stemming from asymmetric information and interest divergence between those who perform tasks (the agents) and those on whose behalf tasks are performed (the principals). A rent can thus be captured by the agents at the expense of their principal.[4] The reason is that a low level of output can be due either to a low exogenous "state of nature" or to some misbehaviour (such as a low effort level or corruption) by the agent. In our model, the ministry of finance (MoF) acts as the principal, providing public funds to line ministries (for example, the ministry of education, the health ministry, or some other public body) to implement a set of actions. The relationship between the MoF and the line ministries is an agency problem subject to asymmetric information both on some external parameters and on the actions performed. In case of low output, the MoF is not in a position to distinguish the cause, unless it uses some form of audit. The principal's problem is to design the contract that most efficiently forces the agent to meet the requirements. The contract must therefore specify a level of output (depending on the state of nature) associated with a certain level of transfer, as well as some control and sanction parameters.

The MoF has a number of instruments at its disposal to limit rent-seeking behaviours. These include, in particular, internal controls within the line ministry. In the so-called francophone system of public expenditure management (prevailing in most of French-speaking Africa), the MoF usually places some of its employees in the line ministries, and their duty is to check that the operations performed by the line ministries comply with the contract. In the anglophone system (prevailing in most anglophone African countries), the approach is different: the line ministry is accountable for its performance *ex post* (and this is verified by the court of audit) and tries to prevent non-compliance by having some of its own employees check the operations of others. In a sense, the head of the line ministry becomes the principal and its employees are the agents. In most cases (anglophone, francophone or other systems), outcomes are verified *ex post* by a court of audit in charge of analysing the performance of the line ministry and reporting its findings to the relevant authorities (usually parliament).[5] If corruption is detected, the official concerned will be punished by disciplinary action or through the judicial system, sometimes entailing a hefty penalty (such as the "*mise en débêt*" in France).[6] These examples suggest that the

literature on incentives and contract design provides the background for a potentially very useful model to analyse PEM issues and guide reforms.[7] However, most traditional models do not exactly fit the realities of PEM systems, and adjustments are needed to take into account their specific features and constraints. These will be identified and elaborated upon as the model is developed.[8]

The article is structured as follows. We interpret a PEM system in light of the assumptions commonly found in the principal-agent literature in Section 2. We then present the basic features of the model in Section 3. Section 4 discusses *ex post* audits and their value in developing countries, while Section 5 introduces *ex ante* controls. We conclude in Section 6.

## 2. Interpretation of PEM under the principal-agent theory

### 2.1. *The contract*

We base our analysis on standard principal-agent models involving supervision (Kofman and Lawarrée, 1993, 1996; Khalil and Lawarrée, 2006). We essentially focus on the control of line ministries or assimilated bodies by the MoF, which is supposed to represent the public interest. Line ministries can be seen as agents of the MoF (the principal) because they are required to produce a certain level of public output – including the quality of this output – in exchange for their budget appropriation. The pair "expenditure programme – budget appropriation" can be interpreted as the two components of the contract between the MoF and the line ministries.[9] The objective of the MoF is to induce the line ministries into implementing their expenditure programmes, while the line ministries pursue their own objectives. That relationship entails both hidden actions (*e.g.* the productive "effort" of the civil servants, possible perquisite consumption, or corruption) and hidden information (*e.g.* the exogenous productivity of that particular sector of the economy), with the agents having the informational advantage over the principal. Hidden information could also refer to poor programme design, which would lead to inefficiency and would be difficult to dissociate from the inefficiency originating from a weak PEM system.[10]

Importantly, the principal-agent model does not allow for a cheating principal and, while it may be argued that the case can occur, we do not consider it in this article.

As already indicated, a number of government operations can be assimilated to principal-agent relationships. For example, one could consider that the minister (who is the head of the ministry, but also a political appointee) heading the line ministry is a principal whose objective is to make sure that his/her agents (the civil servants) implement what he/she has promised to do. One could also consider that the parliament is the principal, whose objective is to make sure that

the government (the executive) implements the government's programme. Yet another example would be to model the central government as the principal, while the subnational governments are the agents. A paper by Ahmad, Tandberg and Zhang (2002) looks at this issue, using a principal-agent framework to analyse incentive structures that best compel local governments to truthfully reveal their ability to implement national programmes. The paper focuses on the optimal contract between both levels of governments. It also insists on the need to have a multi-period game in order to make punishments credible.

An important element of any principal-agent model is to specify an observable that will be the main element of the contract. When the agent is the line ministry, measuring performance should ideally be based on a mix of indicators including output, outcome, and impact. Such information is usually difficult to obtain, and although simply measuring inputs is clearly not satisfactory, they are often the only variable for which adequate data are available. Furthermore, the level of resources in many developing countries, including the HIPCs, is such that broadening the statistics coverage can lead to the undesirable consequence of a serious degradation in the quality of data (*i.e.* information on inputs). Beyond the availability of data, the complexity of performing meaningful measurements, and potential biases linked to the use of performance indicators, it may also not be fully realistic to assume that the MoF has the capability to judge outcomes. Hence, we will hereafter use the term "output" in a general sense, *i.e.* to mean one variable that the MoF is in a position to measure, and this could include outcomes.

It also follows that we formally model a programme budgeting process, although the results also apply to countries that do not explicitly use that approach.[11] In our model, line ministries must make proposals on their priorities, on objectives to be reached, and on corresponding (quantifiable) targets.[12] The line ministry budget proposal is then negotiated with the MoF. However, we do not model the negotiation process. The contract comprises both:

- the required "output" to be produced by the line ministry (in terms of provision of public goods and services) and thus, implicitly, the "effort" required from the line ministry; and
- the line ministry "transfer", *i.e.* its budget appropriation.

A menu of possibilities can be included in the contract to take into account the general economic conditions or make relevant assumptions. For instance, it could be specified that, under a baseline scenario with realistic growth prospects, the line ministries are required to operate with their existing capacities; but, under a more optimistic assumption (the country receives more debt relief, or the economy experiences higher growth), the line ministries could make additional investments. In fact, this is increasingly happening in the context of poverty reduction strategy papers (PRSPs), which

often present a baseline scenario, plus a higher aid scenario based on the resource availability necessary to reach the millennium development goals.[13] We assume that, everything else being equal, the line ministry prefers being granted a large budget appropriation but dislikes the effort associated with the performance requirements.

## 2.2. Agency problems

The agency problem arises from the diverging interests of the MoF and the line ministry and the latter's informational advantage, both on its own actions and on the current state of nature. As standard in the principal-agent literature, the agent's effort is a necessary component of the production function, but entails some disutility. The agent may take unfair advantage of its superior information: if external conditions are favourable, the line ministry could exert little effort and produce a low output, while claiming that this low output is due to unfavourable external conditions. The MoF is not in a position to disentangle the two factors unless it uses some form of audit or supervision. There is thus a risk that the line ministry captures some rent at the expense of the MoF.[14] In the principal-agent literature, this cheating rent generally stems from lowering the level of effort *vis-à-vis* the compensation received. Rents, and possible reductions in public output, compared to what is economically efficient, constitute the agency costs.

In this article, we broaden the interpretation of corruption (generally referred to as the abuse of public office for private benefits) to include misgovernance stemming from the abuse of some information asymmetry. We consider the line ministry's effort in terms of a combination of factors that can be good and bad, including, on the one (good) hand, an efficient and equitable allocation of resources, fiscal transparency measures, and quality of services provided; and, on the other (bad) hand, corruption, consumption of perquisites, mismanagement, and nepotism in the choice of staff or suppliers. This allows us to interpret the cheating rent, not only in terms of reduced disutility from "productive" effort, but also as corruption or misgovernance. For example, if the state of nature is high (say, favourable weather conditions), the line ministry could allocate some resources to unproductive areas or divert monies, if it thinks that the MoF could be led to believe that the state of nature was low. In such cases, rent capture takes place and is possible because of the information asymmetry between the principal and its agent. This interpretation enables us to link our approach with the empirical literature on corruption. Indeed, the latter identifies various factors contributing to corruption, including the overall level of potential benefits from corrupt behaviour, the cost of bribery (including penalties and sanctions), and the bargaining power and extent of discretionary powers of the various actors (Chand and Moene, 1999). Moreover, while cheating (exerting a lower effort) is

probably costless for the agent, we argue in this article that cheating, in the sense of being corrupt, may entail some costs to be concealed. This enables us to make the link with the literature on collusion in organisations (*e.g.* Tirole, 1986) and, in Section 5, we interpret *ex ante* controls by the MoF as increasing the cost of cheating for the agent.

As already stated, the MoF has a number of instruments and strategies at its disposal to limit agency problems. First, it can use incentive schemes, designed solely on observable information, and promise to grant the line ministry a transfer equivalent to the sum of a suitable compensation for the line ministry's effort and an informational rent (which depends on incentive compatibility constraints) in case of high productivity.[15] If such a contract exists, it prevents the line ministry from exerting little effort – but at the expense (for the MoF) of a loss equivalent to the informational rent, in addition to a distortion created by requiring a lower level of effort in some occurrences of the productivity factor. Although commonly applied to models of the corporate world (*e.g.* granting board members bonuses or shares), this strategy is not always directly applicable in the public sector.[16] Alternatively, the MoF can supervise the line ministry using a number of instruments and can threaten it with appropriate sanctions if cheating is detected. The design of the appropriate control system must take a number of factors into account, for instance the choice between *ex ante* and *ex post* (or internal and external) controls, the type of variables to be monitored (input *versus* result indicators), and the choice between systematic or random audits. In our model, there are two unobservable variables (effort and state of nature). Supervision could thus turn either to the exogenous productivity factor, from the observation of which the agent's behaviour could be inferred (this relates, for instance, to public sector reforms aimed at improving the economic statistical data collection, or to audits targeted at assessing the programme design), or directly to the agent's effort. In this article, we assume that the MoF will audit the line ministry's effort.[17] The timing of the game is the same as in all principal-agent models (see Leruth and Paul, 2006, for more details).

# 3. The basic model

In this article, we will not develop the details of the model but concentrate instead on its main elements and results (for more details, refer to Leruth and Paul, 2006).

## 3.1. *Main assumptions*

The model developed here is close to the literature on supervision, as in Kofman and Lawarrée (1993) and Khalil and Lawarrée (2006). However, it differs in some assumptions so as to better reflect the features of PEM systems. The

added value of this article lies in the practical applications to PEM of the analysis, but also in those differences.

We model the agency relationship between the MoF and one line ministry, both assumed to be risk neutral.[18] The line ministry produces a level of output $x$, which depends on two variables: a random exogenous productivity factor, $\theta$; and the line ministry actions or effort, $e$, such that $x = \alpha(\theta,e)$; with $\alpha_e > 0$; $\alpha_{ee} < 0$; $\alpha\theta > 0$. The realised output is public knowledge, but $e$ and $\theta$ are the line ministry's private information. The external productivity can be either high or low: $\theta_i$ with $i \in \{H,L\}$ and $\Delta\theta = \theta_H - \theta_L > 0$. We also assume $\alpha(\theta_H,e) > \alpha(\theta_L,e)$; and $\alpha_e(\theta_H,e) > \alpha_e(\theta_L,e) > 0$. It is common knowledge that the MoF assigns *ex ante* probability $q$ to the event that $i = H$ (and probability $[1 - q]$ to the event that $i = L$). When state $i$ occurs, the line ministry exerts a certain effort level $e_i$, thus producing an output $x_i = \alpha(\theta_i,e_i)$.

The monetary equivalent of the line ministry's disutility from effort is represented by an increasing and strictly convex function $\psi(e)$, with $\psi_e > 0$ and $\psi_{ee} > 0$. To obtain strictly positive but bounded optimal efforts, we also assume that $\alpha(\theta,0) = 0$; $\psi(0) = 0$; $\lim_{e \to 0} \psi_e(e) = 0$; $\lim_{e \to 0} \alpha_e(\theta,e) = \infty$; $\lim_{e \to \infty} \psi_e(e) = \infty$; and $\lim_{e \to \infty} \alpha_e(\theta,e) = 0$. The line ministry's utility is given by $u = t - \psi(e)$, where $t$ is the transfer (appropriation) it receives, and its reservation utility is normalised at zero. We also assume that $\Delta\theta$ is large enough so that the MoF is always better off in case of high output: $x_H - t_H > x_L - t_L$.

### 3.2. Perfect information benchmark

The MoF's problem is to choose the levels of effort required from, and transfers to be made to, the line ministry for each occurrence of the random factor, so as to maximise the expected output:

$$\underset{e_H, e_L, t_H, t_L}{Max} E(X) = q\left[\alpha(\theta_H,e_H) - t_H\right] + (1-q)\left[\alpha(\theta_L,e_L) - t_L\right] \qquad (P)$$

Subject to the line ministry's individual rationality *(IR)* or participation constraints under each occurrence:

$$t_H - \psi(e_H) \geq 0 \qquad\qquad\qquad\qquad\qquad\qquad IR(H)$$

$$t_L - \psi(e_L) \geq 0 \qquad\qquad\qquad\qquad\qquad\qquad IR(L)$$

Under perfect information, the MoF equates the line ministry's marginal cost of effort with the marginal value of its product: $\alpha_e(\theta_i,e_i^*) = \psi_e(e_i^*)$, with $i \in \{H,L\}$.[19] The transfers are such that both participation constraints are binding: $t_i^* = \psi(e_i^*)$. The MoF can therefore enforce first-best, efficient efforts, and the line ministry gets no rent. Note too that, according to our assumptions, $e_H^* > e_L^*$.[20]

### 3.3. Second-best solutions

Under imperfect information, effort is not observable. It cannot be directly enforced, but must be indirectly induced. Therefore, the MoF must provide the right incentives so that the line ministry produces the highest possible level of effort. The reason is that when one productivity level takes place, the line ministry could cheat by adjusting its effort so as to produce the output corresponding to the other productivity level. We define the effort level $\tilde{e}_L$ such that $\alpha(\theta_H, \tilde{e}_L) = \alpha(\theta_L, e_L)$. This means that, if $i = H$, the line ministry could exert a low effort $\tilde{e}_L$, and thus produce $x_L$, while claiming to receive $t_L$. The "cheating rent" that the line ministry can get in that case is thus equal to $\left[t_L - \psi(\tilde{e}_L)\right] - \left[t_H - \psi(e_H)\right] > 0$. To ensure that it is never optimal for the MoF to shut down the contract in case of low productivity, we also assume $(1-q)\alpha(\theta_L, e_L) > \psi(e_L) - q\psi(\tilde{e}_L)$.

Traditional principal-agent models rely on incentive-compatible schemes to prevent such cheating. The well-known results from this literature apply and can be summarised as follows. At equilibrium, the line ministry does not cheat and gets no rent when $i = L$. In order to induce it into exerting the right effort level when $i = H$, the line ministry must receive an informational rent equal to $\psi(e_L^{SB}) - \psi(\tilde{e}_L^{SB})$ (where superscript "SB" stands for "second best") in addition to its first-best transfer, where $\tilde{e}_L^{SB}$ is such that $\alpha(\theta_H, \tilde{e}_L^{SB}) = \alpha(\theta_L, e_L^{SB})$. While production is efficient when $i = H$, the line ministry underproduces (compared to the full information benchmark) in case of low productivity: $e_L^{SB} < e_L^*$. Requiring a lower level of effort when $i = L$ enables the principal to decrease the informational rent granted to the LM when $i = H$. This reflects the "rent extraction – economic efficiency" trade-off which characterises adverse selection problems. The MoF therefore incurs an agency cost equal to the difference in expected output between the first-best and the second-best solutions and caused by information asymmetry.

### 3.4. Introducing supervision

In order to avoid forgoing the informational rent, the principal can hire a supervisor and reduce the information asymmetry. Usually, this is combined with the threat of penalty if cheating is detected. In introducing supervision, we will, in some respects, diverge from the existing principal-agent literature, so as to better reflect PEM concerns.

In the context of PEM, supervision may take various forms. One can distinguish internal controls (*e.g.* MoF or line ministry agents responsible for ensuring that expenditures and procurement are performed according to the rules) and external controls (*e.g.* a court of audit reporting to parliament). Controls may take place *ex ante* (*e.g.* comptrollers issuing visas to allow expenditure, or automatic safeguards preventing line ministries from exceeding

budget appropriations)[21] or *ex post* (*e.g.* auditors checking the reliability of fiscal data or the performance of public spending). Different types of controls can be combined. For instance, going back to the example of PEM systems in Africa, it is worth noting that the so-called francophone system rests on the principle of separation between the person who initiates spending (the *ordonnateur*) and the person who pays it (the accountant or *comptable*). The system relies on centralised *ex ante* controls from the MoF, which take place at various stages of the expenditure process and mainly focus on the conformity of spending with regard to procedures and budget appropriation. Anglophone countries, on the other hand, have inherited a decentralised management system, where the line ministries' accounting officers are responsible for budget execution. *Ex ante* expenditure control is mainly exercised by the issue of periodic warrants by the MoF (cash management). The anglophone system relies on independent *ex post* controls by an auditor-general. In practice, however, notwithstanding those conceptual differences and institutional arrangements, both systems have proven to perform poorly (Bouley, Fournel and Leruth, 2002; Moussa, 2004; and Lienert, 2003).

In the next two sections, we introduce two types of supervision. We first study the case of *ex post* audits (Section 4) and explain how a standard principal-agent model with audit may be of interest for the design of PEM systems. We then move to *ex ante* controls (Section 5). In doing so, we assume that *ex ante* controls increase the cost of cheating for the line ministry, as is done in the literature on collusion.

## 4. *Ex post* audits

In this section, we introduce an *ex post* auditor, costing z whether or not it identifies cheating.[22] It could be an external or an internal auditor, in which case the audit cost may be interpreted as the opportunity cost of using MoF resources for controlling the line ministry instead of doing other tasks. The auditor observes an imperfect signal on the line ministry's effort (for instance, this may be done through a review of accounts to check if there has been corruption). We assume that the signal can take two values: "has complied" or "has cheated". The latter occurs only when the line ministry has indeed cheated, while the signal can report compliance by mistake. The monitoring function is such that σ denotes the probability of detecting actual cheating, in which case the line ministry is imposed an exogenous penalty $P$. With the introduction of supervision, the contract specifies not only the transfers and expected outputs (and thus, implicitly, the expected effort levels), but also the probability of audit. We assume that the MoF commits to audit with probability γ after $x_L$ has been observed. When productivity is high, the line ministry may cheat with probability $m$. Given that output is low, the probability that the line ministry has cheated can be written as $\phi = qm/[(1 - q) + qm]$. We also assume

that the auditor is honest and does not collude with the line ministry.[23] An appealing interpretation is that developing countries are often subject to donors' external auditors, which are supposed to be honest (or not in a position to negotiate with the line ministry).

We discuss three possible regimes that the MoF can implement:[24] *i)* the "cheating-proof" regime, which corresponds to the optimal, incentive-compatible contract when commitment is credible; *ii)* the "cheating-inducing" regime, which is not optimal but, as we argue, matches the situation in some developing countries; and *iii)* a "no commitment" case, in which there is no formal commitment to audit at the time of offering the contract, which results in a mixed strategy equilibrium. In Subsection 4.4, we discuss some applications of the model in terms of public expenditure management and compare these regimes on the basis of their relative costs and benefits.

## 4.1. *The cheating-proof regime*

Traditional principal-agent models with credible commitment use the revelation principle to determine the optimal contract. The revelation principle asserts that, to find the optimal payoff of a problem with asymmetric information, one can, without loss of generality, restrict it to the incentive-compatible, individually rational scheme where every agent truthfully reveals his/her private information. To put it simply, this means that the principal can do no better than offer an incentive-compatible contract, which therefore **deters** cheating. Under these circumstances, the penalty is never imposed at equilibrium, but its existence out of the equilibrium path acts as a deterring threat and prevents cheating.[25]

This approach is only applicable to settings in which the principal is able to credibly commit to any outcome of the contract. In a PEM system, the existence of a court of audit may be viewed as a commitment tool, enabling the MoF to make the credible commitment, at the time of offering the contract, that it will audit the line ministry at the end of the fiscal year with a given probability, which can be either probability one (systematic audit) or below one (random audit).[26] Under this regime, there are conditions where the audit threat is such that it prevents the line ministries from cheating.

Formally, the MoF's problem is to choose the levels of transfers, required efforts and audit probability so as to maximise the expected output:[27]

$$\underset{e_H, e_L, t_H, t_L, \gamma}{Max} \; E\left(X\right) = q\left[\alpha\left(\theta_H, e_H\right) - t_H\right] + \left(1 - q\right)\left[\alpha\left(\theta_L, e_L\right) - t_L - \gamma z\right]$$

Subject to *IR(L)*, *IR(H)*, and the following incentive compatibility *(IC)* constraint:[28]

$$t_H - \psi\left(e_H\right) \geq t_L - \psi\left(\tilde{e}_L\right) - \gamma\sigma P \qquad\qquad IC(H)$$

Note that the term $\gamma\sigma P$ relaxes the *IC(H)* constraint, compared to the second-best case. Hence, provided its cost is not too high, we obtain the intuitive result that audit benefits the principal.

As in Kofman and Lawarrée (1993), we find that the optimal contract exhibits qualitatively different types according to the value of the parameters (although the specific thresholds differ, as we have slightly different assumptions; see Leruth and Paul, 2006, for proof and complete results). These types are characterised by different rents, productive distortions and audit probabilities. The shift from one type to the other rests on a comparison between the expected benefits from audit (the penalty and reduction of rents and distortions) and its cost, with both the former and the latter depending on exogenous (country-specific) parameters. The shadow cost (Lagrange multiplier) of the *IC(H)* constraint is a crucial variable in determining the optimal institutional setting[29] and it is only when the benefits from audit exceed its cost, so that $q\sigma P > (1 - q)z$, that it is profitable to audit. In that case, the optimal probability of audit, rents and productive distortions vary, and three sub-types can be identified:[30]

- "**Rent extraction**" **(RE) scheme:** The MoF audits with probability one, the production distortion when $i = L$ corresponds to the second-best solution, and the line ministry gets a rent when $i = H$ (the rent is equal to that of the second-best solution reduced by the expected penalty).

- "**Effort adjustment**" **(EA) scheme:** The MoF still audits with probability one, but the line ministry no longer gets a rent, and the productive distortion is lower compared to the second-best solution (*i.e.* the effort level lies between the second-best and the first-best solutions).

- "**Random audit**" **(RA):** The MoF decreases the probability of audit (*i.e.* $0 < \gamma < 1$) and still deters cheating, while the productive distortion is smaller than in the second-best solution.[31]

In other words, if it is efficient, audit reduces the agency costs (distortions and rents) associated with the second-best contract, while still deterring cheating. When the expected penalty is relatively low, the MoF must concomitantly use other incentives to prevent cheating. For example, our model suggests that the line ministry could get a rent in case $i = H$, which would then lead to distorting production when $i = L$. Nevertheless, audit enables the MoF to reduce the rent granted to the line ministry compared to the second-best level (hence the label "rent extraction"). When the expected penalty rises, the threat increases for the line ministry so that the MoF can reduce the degree of mobilisation of other incentives. In our model, the line ministry can no longer get a rent, and so the productive distortion can be reduced. Audit thus increases the effort requested from the line ministry in case of low productivity (hence the label "effort adjustment"). Also, when the

expected penalty is large, the MoF can reduce the probability (thus the cost) of auditing (while still deterring cheating). Finally, there remains a production distortion when $i = L$, reflecting the trade-off between efficiency and the cost of audit.

As already mentioned, these results show that, when the expected penalty is relatively small, the MoF must concomitantly use other incentive tools, such as informational rents, in addition to audits to be able to prevent cheating. As the penalty threat increases, other incentives become less necessary. When the penalty is very high, even if the MoF reduces the probability of audit, the line ministry will not find it profitable to try to cheat. In developed countries, like France for instance, this principle is at the root of the "sampling" of expenditures and agencies to be audited.

Finally, note that this model may also be adapted to take into account other constraints faced by the MoF. For example, we could think of a situation where the MoF is obliged to comply with some minimum requirements in terms of output (for example, the "education for all initiative" or the provision of some basic health care package), so that it cannot tolerate that the line ministry produces below $x_L^*$. The MoF would still have to fulfill *IC(H)* at the lowest cost (considering the additional constraint), which could be done through a trade-off between granting the line ministry a higher rent when $i = H$ or raising the audit probability (if it is possible).

## 4.2. *The cheating-inducing regime*

When audit is not optimal, we have seen that the MoF should offer the second-best contract, which is characterised by an informational rent paid to the line ministry when the state of nature is high. However, in reality, there may be circumstances preventing the MoF from offering that contract. For instance, the MoF may not be in a position to grant an informational rent when the budgeting system is input-based (line item) or if it is confronted with a tight cash constraint. Besides, one fundamental difference between private sector operations (which have typically been used to illustrate the principal-agent theory) and public sector operations is often that output is easier to quantify when it is sold on the market. Public sector output often is not easy to quantify and must then be estimated at high cost and with uncertainty (how many children have actually learned to read and write, how many have been vaccinated).

There may also be political pressure or legal constraints forcing the MoF to use *ex post* controls by the court of audit, even if it is not efficient to do so. For instance, nearly all sub-Saharan African countries possess a court of audit. Yet, findings from country financial accountability assessments (CFAAs) show that these institutions often suffer from important weaknesses, ranging from

lack of independence to poor capacity (World Bank, 2004b), which impair their ability to deter cheating in a significant manner. The solution recommended by the principal-agent theory is to grant premiums to line ministries so as to induce them not to cheat. But this solution is hardly (if at all) observed in reality. Rather, one often observes the coexistence of a weak supreme audit institution and high levels of cheating (including corruption), with few positive incentives such as performance premiums.

To look into cases where there exists a court of audit which is not effective in deterring cheating, we now introduce a "cheating-inducing" regime: the MoF uses audit but is not able to deter cheating. We are aware that this regime is not optimal for the principal with respect to the constraints usually considered in the literature, as it could get the same output without incurring the cost of audit. However, we believe that it adequately reflects the situation of some developing countries (see Subsection 4.4 on that issue). Therefore, we consider the case in which the MoF audits, but the penalty and audit probability are such that $t_H - \psi(e_H) < t_L - \psi(\tilde{e}_L) - \gamma\sigma P$. In that case, it is always in the interest of the line ministry to cheat when feasible, and the actual output is always low. As audit takes place, penalties will be imposed at equilibrium, but they will not be a sufficient enough threat to deter cheating. To come back to the model, the cheating-inducing regime takes place when the MoF commits to audit with a probability $\gamma^{CI} < \left[ \psi(e_L^*) - \psi(\tilde{e}_L^*) \right] / \sigma P$ (where the superscript "CI" stands for "cheating-inducing").

Characterising the cheating-inducing regime is interesting in that it helps understand the reasons why the MoF is unable to deter cheating. This happens in the following cases:

- The audit probability $\gamma^{CI}$ is deliberately chosen at a level that is too low.
- The parameters preclude sufficient audit – in particular, if the expected penalty is low compared to the rent, the theoretical $\gamma$ which would help deter cheating may turn out to be higher than unity, which is irrelevant; this case may also occur when the effectiveness of audit ($\sigma$) is too low.
- The MoF does not make use of concomitant incentives (*e.g.* rents and distortions).

Finally, note that if $z > q\sigma P$, audit increases the agency cost. The weight of this component could increase when the MoF supervises several line ministries. For instance, if some ministries are less critical than others, and could reasonably function with low production levels, agency costs could be reduced by offering cheating-inducing contracts, thus saving on auditing costs. The money saved could be used to provide the incentives to the priority line ministries and offer them a cheating-proof contract, thus ensuring high production (when $i = H$) in these sectors.

### 4.3. *The case of "no commitment"*[32]

So far, we have considered that the MoF could credibly commit, at the time of offering the contract, to audit the line ministry with a given probability once output is observed. However, commitment to audit suffers from a time inconsistency problem: the contract determined by the revelation principle is optimal *ex ante* but entails inefficiencies *ex post*. Indeed, as the contract deters cheating, the audit cost must be incurred without any compensation in terms of collected penalty. Moreover, as the contracted effort is not efficient when i = L, one obtains a Pareto improvement by renegotiating the contract once information has been revealed. This time inconsistency reduces the credibility of the commitment to audit, all the more if the MoF is facing a tight budget constraint.[33] Other reasons may also contribute to preclude the commitment to audit, for instance the difficulty of verifying whether the principal did indeed abide by its committed random audit strategy (Mookherjee and Png, 1989) or, in practice, the absence of adequate legal institutions. In particular, we have so far assumed that the existence of a court of audit in the country consisted in a commitment control. Yet, the inefficiency of the court of audit and related institutions may reduce the credibility of the commitment. Nevertheless, in such situations where the MoF cannot credibly commit to auditing, it can call upon some external auditors for specific tasks.

In this subsection, we drop the assumption that the MoF commits to auditing at the time of offering the contract. However, once output is observed, the MoF can decide to audit if it proves to be efficient *ex post*, *i.e.* if the MoF expects to get "value for money" out of the audit. Indeed, when the output is produced, it is too late to deter cheating – but the MoF can still earn a penalty if cheating is detected. The MoF will be willing to audit only if the expected penalty is at least equal to the audit cost. Moreover, the mere expectation that the MoF may audit will reduce the line ministry's incentive to cheat.

Formally, with no commitment, the revelation principle cannot be used and audit must be optimal *ex post* to justify its cost. Hence, cheating may occur in equilibrium (*i.e.* the probability of cheating is positive), and the MoF can expect to collect a penalty. The MoF's problem is to choose the levels of transfers, effort and audit probability so as to maximise its objective function:

$$\underset{e_H, e_L, t_H, t_L, \gamma}{Max} E(X) = q(1-m)\Big[\alpha(\theta_H, e_H) - t_H\Big] + \Big[(1-q) + qm\Big]\Big[\alpha(\theta_L, e_L) - t_L + \gamma(\phi\sigma P - z)\Big]$$

Subject to *IR(L)* and the following constraints:

$$(1-m)\Big[t_H - \psi(e_H)\Big] + m\Big[t_L - \psi(\tilde{e}_L) - \gamma\sigma P\Big] \geq 0 \qquad\qquad IR(H)$$

$$m \in \underset{m'}{\arg\max}(1-m')\Big[t_H - \psi(e_H)\Big] + m'\Big[t_L - \psi(\tilde{e}_L) - \gamma\sigma P\Big] \qquad IC(H)$$

$$\gamma \in \underset{\gamma'}{\arg\max}\Big[\gamma'(\phi\sigma P - z)\Big] \qquad\qquad IC(A)$$

Note that, compared to the previous cases, the MoF's objective function and the constraints now encompass the probability of cheating by the line

414

ministry (Subsections 4.1 and 4.2 correspond to corner solutions of this general problem). The two *IC* constraints consist of indifference conditions and determine the game, which is played simultaneously by the MoF and the line ministry. Indeed, the line ministry is indifferent between cheating and being honest when $t_H - \psi(e_H) = t_L - \psi(\tilde{e}_L) - \gamma\sigma P$. The MoF observes $x_L$, and is indifferent between auditing and not auditing when $\phi\sigma P = z$.

This problem yields a mixed-strategy equilibrium. Under this regime, production is efficient and the line ministry gets no rent. The equilibrium is obtained when the MoF audits with probability $\gamma = \left[\psi(e_L^*) - \psi(\tilde{e}_L^*)\right]/\sigma P$ and the line ministry cheats with probability $m = [z(1-q)/(\sigma P - z)q]$. Note that penalties are collected at equilibrium and are exactly offset by the cost of audit. Moreover, as it entails no production distortion, the mixed-strategy equilibrium tends to the first-best when the penalty is very large or audit is free.[34]

To our knowledge, the literature has not attempted to directly compare regimes with and without commitment, because commitment to audit is usually considered desirable as it reduces the *ex ante* cost of inducing truthful reporting (Baron and Besanko, 1984). The limited commitment due to the possibility of renegotiating a contract is typically handled by using the renegotiation-proofness principle. The latter, which is somewhat similar to the revelation principle, says that one can, without loss of generality, restrict the set of possible contracts to the class of contracts that are not renegotiated. Renegotiation-proof constraints are thus added to the set of *IRs* and *ICs* (*e.g.* Bolton, 1990; Dewatripont and Maskin, 1990; but see also Aghion, Dewatripont and Rey, 1994). In practice, commitment may entail some costs, including those linked to the creation of an adequate institutional setting such as the creation of a court of audit. Our framework allows a comparison of the agency costs associated with each regime. For instance, one observes that the agency cost of our mixed-strategy equilibrium consists mainly of the loss of production (when $i = H$ and the agent cheats). However, when external audits are cheap, and/or the expected penalty is large, the cheating probability decreases and the mixed strategy becomes a better option. This could reduce the value of establishing a court of audit if it does not yet exist in the country.

## 4.4. Applying the theoretical framework to PEM systems

The sections above explain the different control regimes the MoF can implement, according (notably) to the audit probability it chooses. The optimality of these regimes depends on the value of exogenous, country-specific parameters such as the level of penalty, the quality of the supervision technique, the cost of audit, and the probability of high productivity. This suggests that the need to base the choice of the control design on a good analysis cannot be overemphasised and limits the applicability of "one-size-fits-all" solutions. In this respect, our model provides an analytical framework that can guide PEM reforms, as it allows

for comparing different institutional designs while taking into account the constraints faced by governments.

Generally speaking, the MoF should choose the audit regime associated with the lowest agency cost. Yet, the regimes vary with the institutional setting. The cheating-proof regime corresponds to a situation where the MoF can credibly commit *ex ante* to auditing with a given probability, for instance if a court of audit exists. This regime can be compared directly with the second-best contract, and the choice will depend on the ratio $q\sigma P/(1 - q)z$. The mixed strategy relies on a different assumption: there is no commitment, but the MoF can choose *ex post* whether to resort to external auditors for specific tasks. To compare these two frameworks, one should add the fixed cost of setting up and running the institutions necessary to allow commitment (*e.g.* the cost of creating a court of audit) to the agency cost of the cheating-proof regime.

In theory all the regimes discussed above could be implemented. In practice, however, additional constraints may restrain the choices available to the MoF, especially in developing countries. On the one hand, several factors may contribute to reducing the ratio $q\sigma P/(1 - q)z$ which conditions the value of audit in those countries. When cheating is detected, the penalty faced by line ministries may be very small or rarely enforced (see, for example, Dia, 1996; Lienert, 2003; Moussa, 2004), all the more if discounted at a high rate (because of the time required for implementation) and if the supervision technology is not performing well (low $\sigma$, for instance due to poor fiscal data; see, for example, Bouley, Fournel and Leruth, 2002); the probability of high output may be low, compared to that expected in industrial countries; and the opportunity cost of audit may be high, considering the scarcity of competent human resources. Therefore, while the MoF can probably deter cheating through *ex post* audits in industrial countries, this may not be so easy in developing countries.

Practical constraints may also restrict the MoF's actions:

- A tight cash constraint and/or the framework of line-item budgeting may limit the availability of informational rents.[35]

- The government may have committed to provide a minimum package of services, which prevents production distortions below a certain level (the level obtained under the optimal contract).

- The MoF may be legally or politically obliged to resort to the court of audit, even if it is not working well.

Finally, if audit is not efficient [due to a low ratio $q\sigma P/(1 - q)z$] and if the MoF cannot enforce the second-best contract (for practical reasons), it can only implement the cheating-inducing regime: the most unfavourable for the MoF, and a regime that is never optimal.

Note that our analytical framework may be extended to evaluate various reforms. One can think of reforms to increase the quality of audit, or to use signalling to help determine productivity (*e.g.* collecting economic information on the sector and/or auditing the quality of the approved programme design).

# 5. *Ex ante* controls[36]

We have so far considered that, notwithstanding the risk of being caught *ex post*, cheating is costless for the agent. This assumption is common in the literature because the cheating rent generally consists of a reduction in the effort made by the agent, which remains his/her private information. But the line ministry's effort may also comprise some negative actions (such as corruption), and this leads us to assume that cheating entails some costs to be concealed. This allows us to make the link with the economic literature on collusion in organisations.[37] The literature distinguishes two types of collusion costs, according to whether they are exogenous (*e.g.* negotiation costs, "physical" strategies to divert monies from their intended purposes) or endogenous (*e.g.* costs stemming from the risk of future detection; see Faure-Grimaud, Laffont and Martimort, 1999; Khalil and Lawarrée, 2006). Ways in which the principal can avoid corruption include: *i)* create incentive payments; *ii)* decrease the stake of collusion; and *iii)* increase the transaction cost of collusion (Laffont and Rochet, 1997). In this section, we introduce an exogenous cost of cheating and explain how it affects the constraints of the MoF's problem. In a second step, we interpret *ex ante* controls, undertaken by the MoF before the commitment and/or the payment of the line ministry's expenditures, as increasing the cost of cheating. We then discuss the relative value of *ex post* and *ex ante* controls.

## 5.1. *Exogenous cost of cheating*

We assume that the line ministry incurs a certain cost $\eta \geq 0$ when it cheats. That cost decreases the expected benefits from cheating. If we first consider a model without *ex post* audit, the MoF's problem can be written as:

$$\underset{e_H, e_L, t_H, t_L}{Max}\ E(X) = q(1-m)\left[\alpha(\theta_H, e_H) - t_H\right] + (1-q+qm)\left[\alpha(\theta_L, e_L) - t_L\right]$$

Subject to *IR(L)* and:

$$(1-m)\left[t_H - \psi(e_H)\right] + m\left[t_L - \psi(\tilde{e}_L) - \eta\right] \geq 0 \qquad\qquad IR(H)$$

$$m \in \arg\max_{m'}(1-m')\left[t_H - \psi(e_H)\right] + m'\left[t_L - \psi(\tilde{e}_L) - \eta\right] \qquad\qquad IC(H)$$

We observe that *IC(H)* is relaxed by the cost of cheating. If $t_H - \psi(e_H) < t_L - \psi(\tilde{e}_L) - \eta$, the line ministry will always cheat ($m = 1$) unless appropriate incentives are provided. For instance, the second-best contract in this case would also entail a rent when $i = H$ and a productive distortion when $i = L$, but these would be reduced by the cheating cost. If the cost of cheating is high enough and $t_H - \psi(e_H) > t_L - \psi(\tilde{e}_L) - \eta$, the line ministry will not cheat

($m = 0$) and the MoF will reach the first-best solution [no rent and efficient production, as $IC(H)$ is not binding]. It is thus in the MoF's interest to increase the cost of cheating. This is discussed further in the next subsection.

## 5.2. Ex ante *controls as increasing the cost of cheating*

Most traditional principal-agent models consider monitoring at the *ex post* stage.[38] Deterring cheating (or collusion) *ex ante* is done by granting rewards and/or through a threat of future punishment. For their part, PEM systems also include a series of controls designed to prevent agents from cheating *ex ante*. For example, automatic tools, such as computer-based systems that check the budget appropriations before allowing spending, are designed for that purpose, and a similar role is played by the financial comptrollers placed by the MoF within line ministries. Such controls are particularly extensive in the francophone treasury system (Bouley, Fournel and Leruth, 2002) but also exist in the anglophone system (Diamond, 2002). As observation suggests, however, these control techniques are not perfect, partly because agents are very active in trying to bypass them (Lienert, 2003; Moussa, 2004).

We hereafter interpret setting up *ex ante* controls implemented by the MoF as increasing the line ministry's cost of cheating.[39] We endogenise the cost of cheating as a decision variable of the MoF and first consider a model without *ex post* audit.

Assume that $c$ represents the cost of the *ex ante* controls. It may be interpreted as the cost of implementing and running controls, but also as the economic cost (sometimes heavy) of procedures that may complicate the expenditure process and reduce the line ministry's absorptive capacity.[40] The line ministry's cost of cheating $\eta(c)$ is now endogenous and depends on the controls implemented by the MoF. We assume $\eta_c > 0$, $\eta_{cc} < 0$, $\eta(0) = 0$, and $\lim_{c \to \infty} \eta(c) = \infty$.[41] We limit our analysis to the specification of incentive-compatible schemes (thus, where the MoF deters cheating). The MoF will decide on the levels of *ex ante* control, transfer and effort so as to maximise its expected output, as follows:

$$\underset{e_H, e_L, t_H, t_L, c}{Max} \ E(X) = q\big[\alpha(\theta_H, e_H) - t_H\big] + (1-q)\big[\alpha(\theta_L, e_L) - t_L\big] - c$$

Subject to *IR(L), IR(H)* [which now has the form: $t_H - \psi(e_H) \geq 0$] and

$$t_H - \psi(e_H) \geq t_L - \psi(\tilde{e}_L) - \eta(c) \qquad\qquad\qquad IC(H)$$

Relying on *ex ante* controls is different from *ex post* audits because: i) the performance of *ex ante* controls is "intrinsic" (depending on $\eta$), and does not depend on external factors like the level of penalty; ii) the cost of controls is incurred *ex ante*, whatever the state of nature, while the cost of audit is incurred, if at all, only when a low output is observed; and iii) the decision parameter of the MoF is bounded in the case of *ex post* audits while, in theory,

418

the MoF could increase *ex ante* controls infinitely (although it would not be efficient to do so).

The results and proofs related to this discussion are presented in Leruth and Paul (2006). Once again, the problem remains to respect *IC(H)* at the lowest cost, hence the importance of the shadow cost (Lagrance multiplier) of that constraint. We first show that it is not efficient to increase controls above a certain point where $[1/\eta_c(c)] \geq q$. If the cost of deterring cheating by controls alone is too high (*i.e.* when the minimum level of control that would be necessary to deter cheating, $\eta(c) = t_L - \psi(\tilde{e}_L)$, would be such that $[1/\eta_c(c)] > q$), the MoF must also use other means to deter cheating, including rents and productive distortions. For instance, the line ministry could be granted a rent when *i = H* (the level of this rent decreases with the amount of control by the MoF) and be required to produce the second-best level of effort when *i = L*. In that case, the MoF's expected output corresponds to the second-best solution.

When *ex ante* controls are sufficiently effective to allow the MoF to deter cheating by raising controls (until a point such that $[1/\eta_c(c)] < q$), the line ministry gets no rent. However, it may be profitable for the MoF to distort the required effort, if it can save by reducing the level of controls. The MoF would then choose the level of control and the effort required when *i = L* by comparing their cost, *i.e.* so as to relax *IC(H)* at the lowest cost. The trade-off here is between increased efficiency and the cost of control.

The results do not fundamentally differ from the cheating-proof regime with *ex post* audits. Both types of control are assessed with regard to their ability to deter cheating, and if their cost is too high relative to their benefits, the MoF has to use other means (rents and distortions). Yet, one can compare the relative value of both types of controls on the basis of the shadow cost of each problem's *IC(H)* constraint, *i.e.* the Lagrange multiplier of those constraints at equilibrium. Both indeed measure the difficulty of deterring cheating, and determine the agency costs (cost of control, rents and production distortions) associated with each regime. For instance, it is not profitable for the MoF to increase the level of *ex ante* controls when $[1/\eta_c(c)] > q$, nor to use *ex post* audits when $[(1 - q)z + \lambda\gamma]/\sigma P > q$. An analysis of shadow costs reveals that both thresholds depend on the probability that *i = H*. It also reveals that the higher the probability of high output, the higher the incentive to use controls to guarantee it. Second, the effect of *ex post* audits on relaxing the constraint *IC(H)* is mitigated by the size of the expected penalty. As already mentioned, penalties are low and/or not enforced in developing countries, and *ex post* audits may not be able to deter cheating, making *ex ante* controls more effective. Third, if penalties are a credible threat, *ex post* audits may prove to be effective because their cost is incurred only when the principal observes a low output (*i.e.* with probability $[1 - q]$ under a cheating-proof regime) contrary to *ex ante* controls which are imposed in a nondiscriminatory manner.

### 5.3. *Integrating* ex ante *and* ex post *controls*

We now briefly consider (without explicitly solving) the case in which the MoF can deter cheating through a combination of *ex ante* controls and *ex post* audits. Limiting ourselves to incentive-compatible schemes, the MoF's problem is to maximise the expected output:

$$\underset{e_H, e_L, t_H, t_L, \gamma, c}{Max}\ E(X) = q\left[\alpha\left(\theta_H, e_H\right) - t_H\right] + (1-q)\left[\alpha\left(\theta_L, e_L\right) - t_L - \gamma z\right] - c$$

Subject to *IR(L), IR(H)* and:

$$t_H - \psi(e_H) \geq t_L - \psi(\tilde{e}_L) - \eta(c) - \gamma\sigma P \qquad\qquad IC(H)$$

Under a combined approach, the most interesting case occurs when the MoF simultaneously uses both types of controls. This may only take place when $\lambda = 1/\eta_c(c) = [(1-q)z + \lambda\gamma]/\sigma P \leq q$. This may be interpreted as linking the cost-effectiveness of each type of control and comparing their costs (resp. $c$ and $[(1-q)z + \lambda\gamma]$) to their effectiveness in deterring cheating (which depends on $\eta$ and $\sigma P$). When the MoF uses both audits and controls in combination, it will thus equate the relative contribution of each type of control. Finally, the optimal level of *ex ante* controls and probability of *ex post* audits, as well as the production distortion when $i = L$ and the possible rent when $i = H$, will be determined simultaneously, so as to satisfy the constraint *IC(H)* at the lowest cost.

## 6. Conclusion

We have argued that the principal-agent theory offers a powerful analytical framework to better understand PEM systems and guide their design in developing countries. The model discussed in this article equally applies to "managerial" systems relying on *ex post* audits (in the British tradition) and to systems relying more on *ex ante* controls (in many francophone African countries). It allows for comparisons between institutional settings (*e.g.* depending on whether or not the MoF is able to commit to a certain audit probability) and types of control (*e.g.* comparing the effectiveness of *ex post* audits and *ex ante* controls) by examining the cost-effectiveness of various tools available to the principal to deter cheating. However, this often entails some productive distortions, which result from a trade-off between economic efficiency, on the one hand, and the cost of control and/or an informational rent, on the other hand. Finally, we have interpreted corruption and the lack of governance as "rents" captured by the line ministries at the expense of their principal as a result of the informational advantage. By assuming that the agent's effort encompasses productive effort, as well as negative actions such as those related to corruption, we have linked the model to the literature on collusion in organisations.

The model shows that several regimes can exist and that their optimality depends on country-specific parameters, hence the importance of basing the choice of a PEM system on a detailed analysis. Nevertheless, it is possible to draw a few general lessons that can help PEM advisors address some important issues:

- **Ex post controls** (which we mostly assimilate to a court of audit in this article) should be used up to the point where their marginal cost is equal to their return in terms of improved economic efficiency. This will depend on several parameters such as arbitrarily low or ineffective penalties (often the result of social pressure or a weak judicial system in developing countries). Rather than setting up a court of audit (they do not exist in all countries), it may then be profitable to rely on other tools such as external, private audit firms, which increase the cost of cheating for the line ministries. In certain conditions, we also stress the importance of setting up a court of audit so as to make the commitment assumption credible and, in conjunction with better funding, increase its activities, thereby increasing the deterring aspect of the threat of punishment.

- The effectiveness of **internal** controls is similarly determined by cost-benefit considerations, but money spent on internal controls tends to be more effective than money spent on *ex post* controls in developing countries. An important parameter is the extent to which these controls can be bypassed, for example through the use of extraordinary procedures. The cost of internal controls should be assessed carefully, taking into account not only the cost of additional controllers or systems, but also the economic cost due to the resulting slowdown of the expenditure process.

- In countries where the efficiency of both internal and external controls is dubious, theory recommends that the line ministry should be granted an "**informational rent**" in the form of a transfer above the compensation for the effort made. However, in practice, and beyond the difficulty of implementing such schemes in a public sector context in many countries, the efficiency of informational rents may be reduced if appropriate performance measures on which to base the contract between the MoF and the line ministry are unreliable or even unavailable.

- The model may also help **sequence reforms**, although we do not develop this aspect in the article. As causes for the poor performance of the PEM system are identified, it is possible to decide what measures should take priority. For instance, if the MoF is not in a position to deter cheating by introducing internal controls, nor to grant informational rents, it is trapped in the so-called cheating-inducing regime. A first step could be to announce that private audit firms will be hired. In our model, this would relate to implementing a mixed-strategy equilibrium, which tends to be an easily

implementable and cheaper solution (as it incurs no fixed cost) to reduce the extent of cheating. If the country lacks reliable fiscal data, *ex post* audits are not very effective, and the priority should be to reinforce the accounting system before reinvigorating the court of audit.

Finally, although the principal-agent theory provides very interesting insights for the design of PEM systems, we have only considered a few aspects and many more are worth exploring. For instance, a principal-agent analysis could be applied to the allocation of resources for control purposes between different line ministries (for example because they have different probabilities of cheating). Future research could also focus on the dynamic aspects of PEM design and take into account the repeated interactions between the MoF and line ministries at the time the contract is prepared. Although not easily tractable, the realities of the negotiation process between the MoF and the line ministries are very complex, with some line ministries being better informed than others.

### Notes

1. However, there is a growing literature on performance, programme and output budgeting, which basically aims to improve the information available on the effectiveness of public expenditure, and hence helps improve performance through enhanced accountability.

2. This assessment was initiated in 2002 and recently updated (see IDA/IMF, 2002 and 2005).

3. More recently, several bilateral donors and multilateral institutions have set up the public expenditure and financial accountability (PEFA) programme. It aims to build a strategic and collaborative approach to assessing and reforming partner countries' public expenditure, procurement and financial accountability systems, and identifies a set of performance indicators and benchmarks in order to help address developmental and fiduciary objectives. It has developed a public financial management (PFM) performance measurement system, and assesses PFM systems against six critical objectives: budget realism; comprehensive, policy-based budget; fiscal management; information; control; accountability and transparency. The sixteen criteria are presented in Appendix I of Leruth and Paul (2006).

4. This interpretation is consistent with that of political and social sciences, which refer to the broader notion of "rent capture" rather than corruption.

5. In some countries, such as Belgium and Lebanon, the court of audit may also perform some *ex ante* control. In this article, we restrict our attention to the functions of *ex ante* in contrast to *ex post* controls, irrespective of the institution performing them.

6. The "*mise en débêt*" is a tool to make public accountants personally responsible for financial wrongdoings discovered in their management of public funds by the court of audit or similar body.

7. Note that assuming a strict agency relationship between the MoF and the line ministry is a simplification of realities. A powerful line ministry could often play a

very important role in budget negotiations, simply because it has more knowledge about requirements in its own area than the MoF. This could be particularly true at the time of budget preparation, where the line ministry could consider itself as an equal partner to the MoF.

8. Potter and Diamond (1999) and Bouley, Fournel and Leruth (2002) discuss several of these constraints.

9. This is in line with the approach adopted, for instance, in the Australian budgeting system, where so-called service and resource allocation agreements are prepared and implemented (New South Wales Treasury, 2000).

10. For example, if the health system in a country does not perform well, say in terms of vaccination ratios, it can be because the health ministry focuses on other things (and could do those efficiently). It could also be because the money appropriated for the purchase of vaccines gets "lost" in the system. A weak PEM system would generally refer to the latter. See, for example, Gupta and Verhoeven (2001).

11. To quote a very practical definition from the New South Wales Treasury: "Output budgeting involves the Executive Government explicitly 'purchasing' outputs from program and service delivery agencies (the 'providers') in order to achieve desired Government outcomes […]. With performance budgeting, the Executive Government funds (or 'purchases') an agency's program and delivery plan (a set of program and delivery strategies) which the agency has developed in order to achieved desired Government outcomes" (New South Wales Treasury, 2000, p. 13).

12. Ideally, these objectives are defined in the context of a medium-term (three-year) framework and based on a comprehensive macroeconomic model. A multi-year budget framework has the potential to improve incentives, for instance by allowing the introduction of intertemporal competition across agents (Ahmad and Martinez, 2004). Although we will not address this issue in the context of the article, it is important to note that the lack of a proper framework for medium-term budgeting has also been identified by the IMF and the World Bank as an area that requires substantial strengthening.

13. The PRSP approach was launched in 1999 in the context of the HIPC initiative. A PRSP describes a country's macroeconomic, structural and social policies and programmes to promote growth and reduce poverty, as well as associated external financing needs. Its preparation and implementation now often condition the release of aid funds and debt relief.

14. Hereafter we use the term "**informational** rent" when referring to the supplementary premium that the line ministry is deliberately granted as an incentive to exert high effort. We use the term "**cheating** rent" to refer to the amount illegally diverted, notably through corruption.

15. In this article, we use indifferently "performance payment" (contingent on observable/verifiable results) and "informational rent" although the latter is, in principle, more general (the difference between the expected utility of an agent with private information and his/her reservation utility).

16. Incentive schemes may be used in public companies (see, for instance, the regulation theory following Laffont and Tirole, 1993) and also, to some extent, in customs administrations (on the theoretical side, see, for instance, Besley and McLaren, 1993).

17. A model close to ours, which combines adverse selection and moral hazard, predicts that monitoring the agent's action is strictly preferable to auditing private information (Kessler, 2000).

18. In practice, the MoF would try to maximise the joint output of several line ministries. By restricting the model to one line ministry, we assume that the MoF treats all line ministries equally. The case of several line ministries is indirectly handled when the probability of audit is below one (notably in the mixed strategy equilibrium), which may be interpreted as follows: the MoF can possibly audit only a certain number of line ministries, and each line ministry chooses its cheating level considering that probability of being audited.

19. Superscript "∗" stands for first-best values.

20. This could be interpreted as an absorptive capacity constraint: when the state of nature is high, the line ministry must work harder to absorb a larger appropriation.

21. For example, if a government is not allowed to use an overdraft facility with the Central Bank, the rule can be interpreted as a safeguard.

22. Assuming the auditor is paid only when reporting cheating is not relevant in a PEM system, although there are many instances where a bonus is given when cheating is detected (for example, customs employees detecting a fraud have a right to a portion of the tax recovered in many countries).

23. Note that this assumption is also made by Kofman and Lawarrée (1993) in their regime with one costly truthful auditor, as well as in Khalil (1997), for example.

24. We use a terminology similar to that of Eskeland and Thiele (1999) who refer to collusion-proof and collusion-inducing regimes, as we refer to cheating-proof and cheating-inducing regimes.

25. This may seem a little remote from reality. However, the revelation principle characterises the optimal payoffs and can be seen as the "truth-telling map" of a complex mechanism where cheating and punishments occur (Kofman and Lawarrée, 1993, p. 648).

26. If there are several line ministries, the latter situation could correspond to the case where the MoF announces it will audit a certain number of line ministries – so that each line ministry knows, *ex ante*, with which probability it will be audited at the end of the year.

27. The probability of cheating and the expected penalty do not enter the principal's objective function as, under this regime, cheating is deterred. However, the expected penalty appears in the *IC* constraint.

28. To be complete, one should also introduce an *IC(L)* constraint, aimed at preventing the line ministry from producing the high output when productivity is low. This would take the form $t_L - \psi(e_L) \geq t_H - \psi(\tilde{e}_H)$. But, as is common in the literature, that constraint is redundant with the others and is therefore not relevant. Note also that, as we use the revelation principle and deter cheating, the optimal solution exhibits no cheating. The penalty is not collected in equilibrium and hence it does not enter the MoF's objective function (although it is present in the constraints).

29. Note that, from the specification of that Lagrange multiplier, one observes that the higher the cost of audit and/or the probability of low productivity, and the lower the penalty, the harder it is to deter cheating.

30. Our results are consistent with the analysis of Kofman and Lawarrée (1993) with a truthful auditor. There are some slight differences, however.

31. The random audit case may be interpreted in a context with several line ministries, where γ represents the probability, for each line ministry, to be audited in the framework of the general auditing policy of the MoF.

32. Mixed-strategy equilibriums have generally been used in the literature on collusion. The latter suggests that, instead of trying to systematically deter (or induce) collusion, it may be efficient to allow it to some extent (*e.g.* Kofman and Lawarrée, 1996; Khalil, 1997; Khalil and Lawarrée, 2006). This may be the case, for example, if there is a positive probability that the agent and the supervisor are honest, so that it may be too costly to provide incentives to systematically deter collusion. We hereafter apply a similar approach to cheating (corruption). We do not model it; but, in a context with several line ministries, the mixed-strategy equilibrium could also yield the optimal contract when some line ministries are honest while others are corrupt, because preventing cheating as if all line ministries were corrupt would be too costly.

33. This argument holds in one-period games. However, in the context of repeated relationships, it is probably in the MoF's interest not to backtrack on its promise to audit with the announced probability, in order to preserve its reputation.

34. Our results are quite intuitive. Both the probability of audit and the probability of cheating are decreasing functions of the expected penalty. The more cheating rent the line ministry can capture, the more the MoF audits. The more expensive the audit, the more the line ministry cheats. The cheating probability is also influenced by the relative probabilities of the productivity occurrences: to keep the MoF indifferent between auditing and not auditing, cheating is increasing with the probability of low productivity. The agency cost decreases with the expected penalty and increases with the probability of low productivity. It is also higher when the difference of production between $i = H$ and $i = L$ is higher.

35. Informational rents could be envisaged in a system of performance budgeting, as they consist of rewarding the line ministry for good performance (above the compensation of its effort).

36. This section deals with an issue not often discussed in principal-agent papers where the focus tends to be on controls run *ex post*.

37. Following Tirole (1986), that branch of the literature studies the potential for side-contracting between a privately informed, cheating agent and the supervisor hired by the principal to control him or her.

38. As an exception to this statement, Strausz (2006) compares the value of monitoring *versus* auditing – but our analysis differs from Strausz's in that we do not assume that *ex post* audit and *ex ante* controls rely on the same technology.

39. The literature on collusion adopts a similar approach when it acknowledges that hiring a collusive auditor is still useful, because it makes shirking costly for the agent, as he/she will have to pay a bribe to falsify the report (*e.g.* Kofman and Lawarrée, 1996).

40. For instance, in Senegal, the procedures for disbursing the Health Decentralization Fund are such that it takes on average ten months for the resources to be at the disposal of the providers. This leaves only two months to the facility to absorb those resources (World Bank, 2004a).

41. The more effective the controls, the higher $\eta(c)$ for any $c > 0$.

## References

Aghion, P., M. Dewatripont and P. Rey (1994), "Renegotiation Design with Unverifiable Information", *Econometrica*, Vol. 62, No. 2, pp. 257-282.

Ahmad, Ehtisham and Leonardo Martinez (2004), "On the Design of Targeted Expenditure Programs", IMF Working Paper 04/220, International Monetary Fund, Washington DC.

Ahmad, Ehtisham, Eivind Tandberg and Ping Zhang (2002), "On National and Supranational Objectives: Improving the Effectiveness of Targeted Expenditure Programs", IMF Working Paper 02/209, International Monetary Fund, Washington DC.

Baiman, Stanley, John H. Evans and Nandu J. Nagarajan (1991), "Collusion in Auditing", *Journal of Accounting Research*, Vol. 29, pp. 1-18.

Baron, David and David Besanko (1984), "Regulation, Asymmetric Information, and Auditing", *RAND Journal of Economics*, Vol. 15, No. 4, pp. 447-470.

Baron, David and Roger Myerson (1982), "Regulating a Monopolist with Unknown Costs", *Econometrica*, Vol. 50, No. 4, pp. 911-930.

Bayart, Jean-François (1993), *The State in Africa: The Politics of the Belly*, Longman, New York, United States.

Besley, Timothy and John McLaren (1993), "Taxes and Bribery: The Role of Wage Incentives", *Economic Journal*, Vol. 103, No. 416, pp. 119-141.

Bolton, Patrick (1990), "Renegotiation and the Dynamics of Contract Design", *European Economic Review*, Vol. 34, No. 2-3, pp. 303-310.

Bouley, Dominique, Jerôme Fournel and Luc Leruth (2002), "How Do Treasury Systems Operate in Sub-Saharan Francophone Africa?", *OECD Journal on Budgeting*, Vol. 2, No. 4, pp. 49-84.

Chand, Sheetal K. and Karl O. Moene (1999), "Controlling Fiscal Corruption", *World Development*, Vol. 27, No. 7, pp. 1129-1140.

Dabla-Norris, Era and Elisabeth Paul (2006), "What Transparency Can Do When Incentives Fail: An Analysis of Rent Capture", IMF Working Paper 06/146, International Monetary Fund, Washington DC.

Dewatripont, Mathias and Eric Maskin (1990), "Contract Renegotiation in Models of Asymmetric Information", *European Economic Review*, Vol. 34, pp. 311-321.

Dia, Mamadou (1993), "A Governance Approach to Civil Service Reform in Sub-Saharan Africa", World Bank Technical Paper No. 225, World Bank, Washington DC.

Dia, Mamadou (1996), *Africa's Management in the 1990s and Beyond: Reconciling Indigenous and Transplanted Institutions*, World Bank, Washington DC.

Diamond, Jack (2002), "The Role of Internal Audit in Government Financial Management: An International Perspective", IMF Working Paper 02/94, International Monetary Fund, Washington DC.

Diamond, Jack (2006), "Budget System Reform in Emerging Economies: The Challenges and the Reform Agenda", Occasional Paper No. 245, International Monetary Fund, Washington DC.

Eskeland, Gunnar S. and Henrik Thiele (1999), "Corruption Under Moral Hazard", World Bank Policy Research Working Paper No. 2204, World Bank, Washington DC.

Faure-Grimaud, Antoine, Jean-Jacques Laffont and David Martimort (1999), "The Endogenous Transaction Costs of Delegated Auditing", *European Economic Review*, Vol. 43, No. 4-6, pp. 1039-1048.

Gupta, Sanjeev and Marijn Verhoeven (2001), "The Efficiency of Government Expenditure: Experiences from Africa", *Journal of Policy Modeling*, Vol. 23, No. 4, pp. 433-467.

IDA (International Development Association) and IMF (International Monetary Fund) (2002), "Actions to Strengthen the Tracking of Poverty-Reducing Public Spending in Heavily Indebted Poor Countries (HIPCs)", March, Washington DC.

IDA/IMF (2005), "Update on the Assessments and Implementation of Action Plans to Strengthen Capacity of HIPCs to Track Poverty-Reducing Public Spending", April, Washington DC.

Kessler, Anke (2000), "On Monitoring and Collusion in Hierarchies", *Journal of Economic Theory*, Vol. 91, No. 2, pp. 280-291.

Khalil, Fahad (1997), "Auditing Without Commitment", *RAND Journal of Economics*, Vol. 28, No. 4, pp. 629-640.

Khalil, Fahad and Jacques Lawarrée (2006), "Incentives for Corruptible Auditors in the Absence of Commitment", *Journal of Industrial Economics*, Vol. 54, No. 2, pp. 269-291.

Kofman, Fred and Jacques Lawarrée (1993), "Collusion in Hierarchical Agency", *Econometrica*, Vol. 61, No. 3, pp. 629-656.

Kofman, Fred and Jacques Lawarrée (1996), "On the Optimality of Allowing Collusion", *Journal of Public Economics*, Vol. 61, No. 3, pp. 383-407.

Laffont, Jean-Jacques and David Martimort (2002), *The Theory of Incentives: The Principal-Agent Model*, Princeton University Press, Princeton, New Jersey, United States.

Laffont, Jean-Jacques and Jean-Charles Rochet (1997), "Collusion in Organizations", *Scandinavian Journal of Economics*, Vol. 99, No. 4, pp. 485-495.

Laffont, Jean-Jacques and Jean Tirole (1993), *A Theory of Incentives in Procurement and Regulation*, MIT Press, Cambridge, Massachusetts, United States.

Leruth, Luc and Elisabeth Paul (2006), "A Principal-Agent Theory Approach to Public Expenditure Management Systems in Developing Countries", IMF Working Paper 06/204, International Monetary Fund, Washington DC.

Lienert, Ian C. (2003), "A Comparison Between Two Public Expenditure Management Systems in Africa", *OECD Journal on Budgeting*, Vol. 3, No. 3, pp. 35-66.

Mbaku, John Mukum (ed.) (1998), *Corruption and the Crisis of Institutional Reforms in Africa*, African Studies Series, Vol. 47, Edwin Mellen Press, New York, United States.

Mookherjee, Dilip and Ivan Png (1989), "Optimal Auditing, Insurance, and Redistribution", *Quarterly Journal of Economics*, Vol. 104, No. 2, pp. 399-415.

Moussa, Yaya (2004), "Public Expenditure Management in Francophone Africa: A Cross-Country Analysis", IMF Working Paper 04/42, International Monetary Fund, Washington DC.

New South Wales Treasury (2000), *The Financial Management Framework for the General Government Sector*, December, Sydney, New South Wales, Australia, *www.treasury.nsw.gov.au/pubs_by_pol*.

Potter, Barry and Jack Diamond (1999), *Guidelines for Public Expenditure Management*, International Monetary Fund, Washington DC.

Reinikka, Ritva and Jakob Svensson (2004), "Local Capture: Evidence from a Central Government Transfer Program in Uganda", *Quarterly Journal of Economics*, Vol. 119, No. 2, pp. 679-705.

République du Cameroun, Institut National de la Statistique (2004), "Enquête sur le suivi des dépenses publiques et la satisfaction des bénéficiaires dans les secteurs

de l'éducation et de la santé. Phase I : Volet Santé", rapport principal des résultats (version provisoire du 12 mars), Yaoundé, Cameroun.

Schick, Allen (1997), *Modern Budgeting*, OECD, Paris.

Schick, Allen (1998), "Why Most Developing Countries Should Not Try New Zealand Reforms", *The World Bank Research Observer*, Vol. 13, No. 1, pp. 123-131.

Strausz, Roland (2006), "Timing of Verification Procedures: Monitoring Versus Auditing", *Journal of Economic Behavior and Organization*, Vol. 59, No. 1, pp. 89-107.

Tanzi, V. (2000), "Rationalizing the Government Budget: Or Why Fiscal Policy Is So Difficult", in Anne O. Krueger (ed.), *Economic Policy Reform: The Second Stage*, University of Chicago Press, Chicago, Illinois, United States, pp. 435-452.

Tirole, Jean (1986), "Hierarchies and Bureaucracies: On the Role of Collusion in Organizations", *Journal of Law, Economics, and Organizations*, Vol. 2, No. 2, pp. 181-214.

World Bank (1997), "Helping Countries Combat Corruption: The Role of the World Bank", Poverty Reduction and Economic Management, World Bank, Washington DC.

World Bank (2004a), "Senegal Public Expenditure Review", Report No. 29357-SN, PREM 4, Africa Region, World Bank, Washington DC.

World Bank (2004b), "Supporting and Strengthening Supreme Audit Institutions: A World Bank Strategy", Financial Management Network, Operations Policy and Country Services, World Bank, Washington DC.

Ye, Xiao and Sudharshan Canagarajah (2002), "Efficiency of Public Expenditure Distribution and Beyond: A Report on Ghana's 2000 Public Expenditure Tracking Survey in the Sectors of Primary Health and Education", Africa Region Working Paper Series No. 31, World Bank, Washington DC.

# An assessment of agency theory as a framework for the government–university relationship

Jussi Kivistö*

*Department of Management Studies, University of Tampere, Tampere, Finland*

The aim of this paper is to use agency theory as the theoretical framework for an examination of the government–university relationship and to assess the main strengths and weaknesses of the theory in this context. Because of its logically consistent framework, agency theory is able to manifest many of the complexities and difficulties that governments face in their attempts to govern universities. Agency theory also offers unique explanations for the government's choice and use of certain governance procedures, low performance by universities and cost growth in the higher education sector.

**Keywords:** agency theory; governance; government; opportunism; university

## Introduction

Agency theory, also known as the principal agent or principal agency theory/model describes the relationship between two or more parties, in which one party, designated as the *principal*, engages another party, designated as the *agent*, to perform some task on the behalf of the principal (Jensen & Meckling, 1976; Moe, 1984; Ross, 1973). The theory assumes that once principals delegate authority to agents, they often have problems controlling them, because agents' goals often differ from their own, and because agents often have better information about their capacity and activities than do principals. Agency theory focuses on the ways principals try to mitigate this control problem by selecting certain types of agents and certain forms of monitoring their actions, and by economic incentives (Kiser, 1999).

Government–university relationships have gone through a period of transition over the past 15 years, especially in Europe. The predominant mode for government governance of universities has seemingly shifted from control and regulation to supervision and enforcement of the universities' self-regulative capabilities (Hölttä, 1995; van Vught, 1997). More specifically, these changes have been characterised by delegation and a shift from hierarchical, authority-based governance structures to contractual, exchange-based governance structures. Increased autonomy and self-regulative capabilities of universities have been accompanied by increases in their accountability to governments (Gornitzka et al., 2004; Trow, 1996). Given this development, one needs to ask two simple interconnected questions:

1. Why do governments need to verify the accountability of universities?
2. Why can't governments just trust universities?

*Email: jussi.kivisto@uta.fi

As a theory characterised by mistrust, control and compliance, agency theory is able to propose a rather simple and straightforward answer to the questions presented. According to agency theory, governments do not trust universities, simply because universities are likely to behave opportunistically if they are not held accountable for the resources they receive.

The purpose of this paper is to assess the main strengths and weaknesses of agency theory when it is utilised as an analytical framework for examining the government–university relationship.

## Literature review

Although strongly influenced by its background in economics (Alchian & Demsetz, 1972; Jensen & Meckling, 1976; Ross, 1973) and political science (Mitnick, 1975; Moe, 1984; Rose-Ackerman, 1978), agency theory is not and it never has been the exclusive property of one particular scholarly discipline or paradigm. Rather, it has a widely applied theoretical and empirical framework for many different disciplines and approaches.

Over the past two decades, the role of agency theory in the mainstream of higher education research has been nominal. Authors such as Ferris (1991), Geuna (1999), Hölttä (1995), Massy (1996), Whynes (1993) and Williams (1995) all have acknowledged and mentioned the principal–agent relationship, but deeper examination of this relationship as an agency relationship was left aside. However, in addition to the works of Kivistö (2005; 2007) and Lane and Kivistö (2008), a growing number of attempts to apply agency theory within the context of higher education has been made, especially in recent years. Scholars such as Gornitzka et al. (2004), Lane (2005), Liefner (2003) and McLendon (2003) have applied the theory within a higher education context and to government–university relationships. To date, higher education studies applying agency theory have utilised it primarily as a conceptual framework, heuristic tool or as an organising concept that aimed to offer insights related to university governance. Owing to the divergent development of agency theory in different disciplines, application of the theory to higher education governance and policy has been somewhat disjointed as scholars using the same broad theory utilise different assumptions based on their own disciplinary perspective, usually political science or economics (Lane & Kivistö, 2008).

## The framework of agency theory

In order to consider the government–university relationship as an agency relationship, the relationship must contain the following three elements:

1. Tasks that the government delegates to a university (i.e. teaching and research);
2. Resources that the government allocates to a university for accomplishing those tasks; and
3. Government interest in governing the accomplishment of the tasks (Kivistö, 2007).

The parties to the agency relationship; that is, a government and a university, also require more detailed specification. The *government*, as a principal, can be defined differently depending on the chosen context and perspective. In the broadest

sense, government may refer to the whole body of national or regional public institutions performing political or administrative functions. When understood in a narrower sense, government can be viewed as a public bureau or agency, such as a ministry or department (Laking, 2005). Bilateral government–university relationships are more common in Europe, where universities are regulated by relatively strong and unitary central governments. The situation is somewhat different in the United States, where State universities are often surrounded by other actors, such as legislature, the governor, higher education commissioner or a coordinating board (Lane, 2005). A *university*, as an agent, is considered to be either a public or a non-profit higher education institution, which has identifiable legal, economical and/or socio-cultural boundaries separating it from the boundaries of the government organisation. Public universities are funded and owned by public authorities and their legal status is public. Non-profit universities are at least partly publicly funded, but ownership and legal status can be private (i.e. not public).

Agency theory makes two important assumptions concerning the agency relationship. There must be *informational asymmetries* and *goal conflicts* present simultaneously in the agency relationship (e.g. Moe, 1984; Waterman & Meier, 1998). Informational asymmetries can be simply considered as a claim that an agent possesses more or better information about the details of individual tasks assigned to him, his own actions, abilities and preferences (cf. Eggertsson, 1990). Academic work is itself inherently surrounded by high informational asymmetries, starting from its core substance, knowledge (Clark, 1984). Understanding the substance of academic work requires a high level of specified expertise and it is not easily replaceable. Some of the total informational asymmetries are also affected by the structural complexity of university organisation (see, for example, Birnbaum (1988) and Clark (1983)). According to Birnbaum (2001) universities are 'complex, non-linear systems, and their responses to chances in one part can have counterintuitive and surprising effects in another' (p. 194). Furthermore, multiple complexities in production technology of universities are causing informational asymmetries (e.g. Cave et al., 1997; Höltlä, 1995; Johnes & Taylor, 1990).

Goal conflicts refer to a situation in which the principal's and the agent's desires and interests concerning certain ends are in conflict with each other and that they would therefore prefer different courses of action (Milgrom & Roberts, 1992). At the level of official goals, there is likely to be general goal conflict between the cultural arguments (universities) and the utilitarian arguments (governments) (Bleiklie, 1998). Universities have traditionally argued that scholarly and academic activities are emancipating forces in the development of open and tolerant societies. Governments have emphasised that universities, whether they are understood as 'bureaux' or 'corporations', exist mainly to provide society with qualified labour, knowledge and research products that contribute to national economic growth (Bleiklie, 1998; Schmidtlein, 2004). The goal conflict resulting from the clash of cultural and utilitarian goal conceptions can be essentially crystallised by the confrontation between the utilitarian claims for accountability and cultural emphasis on academic freedom and institutional autonomy. From the cultural perspective, it is believed that stronger accountability to government weakens the autonomy of the institutions. In addition to the sincere fears of losing academic freedom and institutional autonomy, covert goals may also play a part in the general university and faculty dislike for the government's accountability demands. In extreme cases,

the overt goals may also serve as sentimental 'hobby-horses' created for the purpose of justifying more self-interested covert goals such as 'pursuing prestige, influence, revenue and leisure' (Kivistö, 2007).

### The agency problem

Taken together, the informational asymmetries and goal conflicts constitute the *agency problem* – the possibility of opportunistic behaviour on the agent's part that works against the welfare of the principal (Barney & Hesterly, 1996). The agency problem (also called 'moral hazard problem') may arise in situations in which the principal cannot directly observe the agent's actions and when the self-interested agent pursues his private goals at the expense of the principal's goals (e.g. Barney & Ouchi, 1986; Milgrom & Roberts, 1992). Self-interest can also make the agent reluctant to share performance information with the principal or to motivate the agent to send wrong information to the principal (Bergen, Dutta & Walker, 1992).

Basically, every action and effort of a university could turn into a form of opportunistic behaviour when the true reasons behind this behaviour are self-interested and not in the best interests of the government. This observation includes the behaviour that takes place both at the individual and institutional levels. According to Kivistö (2007), possible manifestations of opportunistic behaviour could include: (1) *shirking by individuals* (e.g. Gomez-Mejia & Balkin, 1992; Ortmann & Squire, 2000; Whicker, 1997); (2) *opportunistic pursuit of prestige* (see Brewer, Gates & Goldman, 2002; Garvin, 1980) and *opportunistic pursuit of revenues* (see Bowen 1970; 1980; Levin & Koski, 1998; Vedder, 2004); (3) *opportunistic cross-subsidisation* (see James, 1986; 1990); and (4) *the distortion of monitoring information* (Bohte & Meier, 2000). All of these reduce efficiency and effectiveness in the use of government resources. For instance, shirking behaviour and opportunistic cross-subsidisation both have the effect of lowering the university's expected output because they employ productive resources for other purposes. University opportunism is also likely to reduce effectiveness, including the quality of research and teaching processes and outputs. Or the shirking activity of faculty members that lead to a constant absence from scheduled instructional tasks may lead to lower learning outcomes and might prolong students' graduation times unnecessarily. Effectiveness is also lost when funds earmarked for expenditure on undergraduate education are opportunistically transferred to subsidise research or other prestige-generating activities (Kivistö, 2007).

### Governing the agency problem

According to agency theory, the principal has two basic options in seeking to control the agent in terms of the contracts to be agreed upon; that is, behaviour-based contracts and outcome-based contracts (Eisenhardt, 1989). When choosing behaviour-based contracts the principal chooses to monitor the agent's behaviour (actions) and then reward that behaviour. The other option, outcome-based contracts, compensates agents for achieving certain outcomes (outputs). Governments perform numerous controlling and monitoring procedures in their relationships with universities, and many of these procedures have a logical analogy with behaviour-based contracts. Here, all of these arrangements are referred as *behaviour-based governance* procedures. In general, such procedures include all those

reporting requests, site visits, reviews and evaluations that focus on monitoring the productive activities, with the primary purpose of informing the government about how universities are 'behaving' in economic and operational terms. As in behaviour-based contracts, the amount of government funding has a connection with the observed behaviour. Therefore, different forms of input-based funding arrangements (i.e. line-item budgeting or input-based formula funding) applied by the governments represent one type of behaviour-based governance procedure. Input-based funding is a form of government funding procedure, whereby the amount of funding allocated is based on different input elements or a university's production processes; that is, indicators that refer to the resources used or the activities carried out by the universities (Jongbloed & Vossensteyn, 2001).

The other option for the government to prevent the agency problem from occurring is to offer output-related incentives to universities. Similarly with outcome-based contracts, the general objective of *output-based governance* is to reduce goal conflicts by aligning the goals of universities with the ones of the government. It is usually organised through performance-based funding practices that are constructed on some output-based funding formula. Because of the intangible nature of teaching and research outputs, governments have been forced to create surrogate measures and proxies, indicators to describe and represent the outputs (Cave et al., 1997). Output indicators derived from teaching activities can include the number of study credits obtained, the number of exams passed, the number of undergraduate and graduate degrees granted and graduates' employment rates. Output indicators derived from research activities can be the number of research publications, research income, the number of patents and licenses received, the number of doctoral students and the number of graduate/doctoral degrees granted (Jongbloed & Vossensteyn, 2001). In addition, the government may also use more complex output-based performance indicators, such as 'value added', 'graduation rate', 'graduation time' and output-connected average cost measures (Cave et al., 1997).

## Agency variables and agency costs

The central challenge for the principal is to choose between different behaviour-based and outcome-based contracts. For this challenge, agency theory presents the two inter-related concepts of 'agency costs' and 'agency variables'. *Agency variables* describe the levels of different internal and external conditions connected to the agency relationship that may have implications for agency costs and contract choice. In other words, agency variables are believed to be able to predict the most efficient contracting choice for a given situation. Although the exact number of agency variables has varied within different research settings, at least four variables – outcome measurability, outcome uncertainty, task programmability and goal conflict – can be identified (cf. Eisenhardt, 1989; Kivistö, 2007).

The fact that teaching and research outputs are both, to a large extent, immeasurable (definitional problems of what are the 'right' or 'true' teaching and research outputs) and uncertain (uncontrollable student behaviour, unpredictable nature of research work) indicates that the government should make use of behaviour-based governance procedures. However, the low task programmability of teaching and research activities (unique and non-repeatable tasks that require high

levels of expertise and creativity) suggests that the government should use output-based governance mechanisms, which can bypass the informational asymmetries related to low task programmability. Furthermore, high or even moderate goal conflicts between the government and university can create serious incentive problems, increasing the possibility of opportunistic behaviour by universities. This would lead to the suggestion of the use of more powerful economic incentives and output-based governance. Because the government seems to face problems with both types of governance procedures, it has to choose between two less than perfect options; either it suffers from problems related to behaviour-based governance or it suffers from problems related to output-based governance. It is a matter of concrete context, opinion and political debate which of these problems should be considered to be the lesser evil.

The concept of *agency costs* has been afforded a range of definitions but, in the broadest sense, they can be understood to be the total costs of different contracting choices plus the costs resulting from agent opportunism. The costs of agent opportunism are the loss borne by the principal caused by the agent acting in his own interest at the expense of the principal (e.g. Tosi & Gomez-Mejia, 1989). Agency costs for government can be defined as the total sum of the costs resulting from governing universities plus the costs incurred because of the opportunistic behaviour of the universities. The total governance costs include the direct and indirect costs associated with the governing procedures. Unfortunately, the costs of governance in a given concrete situation are often impossible to calculate. It is unlikely that government cost calculation systems would be able to count all the costs that are related to the use of a certain type of governance procedure. Nevertheless, these costs can be estimated indirectly and perceived in other than monetary terms. For instance, the cost of governance procedures could be evaluated indirectly as the amount of planning required to establish and to operate them and the number of new staff required, or the new hierarchies their application creates and the observable or estimated dysfunctions these inflict on universities' production behaviour (Jensen & Meckling, 1976; Kivistö, 2007). Owing to the invisible and unperceivable nature of opportunistic behaviour, the costs of detected and undetected opportunism, 'opportunism costs', are even more difficult to calculate, although analytically they are possible to distinguish (Vining & Globerman, 1999). Nevertheless, as a theoretical concept, they could offer interesting perspectives in speculating on the meaningfulness and effects of the government governance of universities.

### An assessment of agency theory

The question of what constitutes 'strong' or 'weak' theory is not simple (see, for example, Kezar 2006). There is no general agreement as to whether the strength of a theory depends more on how 'plausible', 'interesting' or 'aesthetically pleasing' the theory is (e.g. Weick, 1989), or how well it offers 'explanations' and 'predictions' in order to diminish the complexity of the empirical world (e.g. Bacharach, 1989). This is understandable as scholarly preferences and perspectives, as well as the significance and emphasis put on each of the presented evaluation criteria, seem to vary between theories, theorists, paradigms and disciplines. Nevertheless, a relatively neutral position would be to argue that a good theory is a theory that fulfils its

purpose. As the purpose of theories varies so to will the mode of validating the theory (Dickoff & James, 1992; Kezar, 2006).

### The strengths of agency theory

Government–university relationships contain both of the essential conditions that should be present in an agency relationship: informational asymmetries and goal conflicts. Both of these conditions seem to be relevant to an examination of this relationship and they can be operationalised in the context of government–university relationships. The existence of both informational asymmetries and goal conflicts creates favourable conditions for the appearance of the agency problem. The investigation of agency problem can offer a rich variety of insights on issues that are related to established agency relationships and it allows the conceptualising and operationalising of possible forms of university opportunism. In addition, the agency theory also can cause attention to be focused on the productivity effects of opportunistic behaviour by offering alternative explanations for lower levels of performance by universities. By accepting opportunism as one possible explanation for lower university performance, the theoretical perspective for examining issues like accountability, cost growth, efficiency, effectiveness, performance measurement, funding models or quality assessment can become wider, both theoretically and empirically (Kivistö, 2007).

In addition to analysing opportunism by universities, agency theory allows for the categorisation of funding methods, performance measurement instruments, and other monitoring and assessment practices into two mutually exclusive categories (i.e. behaviour-based and output-based governance procedures). The special importance of this categorisation is that it is able to create conceptual links between different governance procedures and the conditions that cause agency problem; that is, informational asymmetries and goal conflicts. These insights allow more systematic and theoretical analysis of the effects of particular governance methods.

As an important part of discerning the governance of agency problem, agency theory introduces the concepts of agency cost and agency variables that are able to offer insights concerning the costs, efficiency and effectiveness of particular governance methods. When choosing between different governance procedures, the government can analyse and make predictions about the applicability and cost of each procedure in light of the agency variables. In addition to their predictive capabilities, the use of these variables offers help both for conceptualising and analysing many of the strengths and weaknesses that are inherent in using particular behaviour- and output-based governance procedures. (For a more detailed discussion, see Kivistö (2007).) Determining the agency costs can include interesting speculation as to whether the costs of governance could, in some cases, exceed the costs of the actual opportunism. If the costs resulting from opportunism remain lower than the governing costs, the best solution for government could be to reduce its governing efforts.

As a whole, all the strengths offered by agency theory to examine the government–university relationship lie in the theory's unique perspective on issues that other theories do not contribute to, to the same extent. Basically, all of the insights that agency theory can offer are related to the question of universities' compliance with the government's goals in exchange for the resources they receive.

Analysing different forms of university opportunism and issues related to governing the opportunism are clearly the strongest and most unique insights that the theory can offer. Agency theory can provide alternative insights by examining the economic characteristics of universities with respect to the behavioural implications for government governance and resource allocation mechanisms. As such, agency theory seems to be able to offer a broad but logically consistent framework about the government's governance of universities, in which different theoretical concepts and approaches can be integrated in a meaningful way.

It is also worth noting that agency theory is not only a tool that is able to strengthen the governing capabilities of the government; it may also help the faculty, departments and universities to understand and develop those means which can demonstrate that they have accomplished (or at least have attempted to accomplish) their tax-funded tasks. One of the key issues seems to be how to meet the legitimate accountability needs of government while safeguarding the institutional autonomy of universities and the academic freedom of the faculty (Kells, 1994). In fact, rather than weakening the role of the academic community, agency theory can, in fact, strengthen the role of universities and faculty, if they choose to take into their own hands those initiatives that demonstrate their accountability. In addition, agency theory and agency variables can offer a theoretically sound framework for analysing the actual and potential shortcomings of funding schemes, performance indicators or assessment methods and, therefore, increase the weight of the argument directed against their use (Kivistö, 2007).

### The weaknesses of agency theory

Like any other theory, agency theory has its critics (see, for example, Donaldson, 1990, 1995; Perrow, 1986) and part of this criticism is relevant in the context of higher education (see Kivistö, 2007). Agency theory has been criticised mostly because of the behavioural assumptions it makes concerning human motivation and behaviour. The critics of agency theory argue that the theory presents too narrow a model of human motivation and that it makes unnecessary negative and cynical moral evaluations about people. According to critics, focusing on self-interested and opportunistic behaviour makes it possible to ignore a wider range of human motives, including altruism, trust, respect and intrinsic motivation of an inherently satisfying task. This criticism also has validity when agency theory is utilised for analysing government–university relationships. If universities are considered only as aggregates of self-interested shirkers, a high level of realism, objectivity and tactfulness will, undoubtedly, be lost.

The greatest weaknesses of agency theory are related to the narrowness of its behavioural assumptions and of the focus of the theory. The fact that agency theory focuses only on self-interested and opportunistic human behaviour means that the theory ignores a wider range of human motives. Even though agency theory does not suggest that self interest-based opportunistic behaviour is the only motivator of human beings, the problem is that the theory partly fails to explain the principal's losses by any factor other than agent opportunism. In addition, agency theory's behavioural assumptions can also limit the scope of the theory. Agency theory pays attention to mainly formal and economic aspects of government–university relationships. However, in addition to their economic character, universities are also socio-cultural organisations as regards the norms, incentives and organisational structure on

which their behaviour is based. Another problem is that agency theory examines agency relationships without questioning the legitimacy or justification of the principal's goals. Although this does not reduce the coherence and robustness of the theory, it is a clear limitation of its scope, especially in the context of higher education.

Agency theory also suffers from narrowness of focus in mirroring the complexity and diversity of the empirical reality (excluding some applications by political scientists). It offers only a limited view by focusing on the bilateral interaction between the government and a university in the context that is, in reality, often surrounded by multilateral networks. In the modern world, a growing number of stakeholders can be found from inside and outside the universities. The fact that agency theory is not able to include third parties, stakeholders or competing principals holistically in its analysis is clearly a great weakness. Agency theory is able to examine only one of the many agency relationships at a time and it gives no suggestions about how to proportionate this relationship to other possible agency relationships. This inability to structure and incorporate the existence of multiple principals and stakeholders creates a danger that can reduce the theoretical and empirical usefulness of the theory.

## Conclusion

Within its limitations, agency theory offers an analytical tool that can benefit a range of parties interested in examining different aspects of government–university relationships: higher education researchers, higher education community, and higher education policy makers. Agency theory offers a theoretical framework for analysing why governments are demanding quality assurance mechanisms and why they sometimes choose output-based funding procedures instead of input-based funding procedures.

Accepting opportunism as one possible explanation for performance failures might generate a more open and useful discussion and investigation about the productive and economic behaviour of universities. By offering theoretical understandings and solutions for the phenomena of inefficiency and cost growth, agency theory is able to underline some of the most topical questions asked by the higher education community and policy makers. For higher education research purposes, the theory possesses enough theoretical parsimony in its concepts but, at the same time, these concepts are broad enough to allow for a wide variety of operationalisations to be made. Despite the narrowness of its focus and assumptions, agency theory offers a clear, logically consistent and dynamic framework for examining government–university relationship from the selected perspective.

Investigating the possibility of developing agency theory in line with other theoretical approaches in the same area of interest should be considered also. Especially, studying the role of information, interests and incentives could provide a much more rich understanding of the complex dynamics of government–university relationships. For example, the study of performance measurement, studies of policy instruments, theory-driven program evaluation and implementation theory could benefit both agency theory and these approaches.

## References

Alchian, A. A., & Demsetz, H. (1972). Production, information costs, and economic organization. *The American Economic Review*, *62*(5), 777–795.

Bacharach, S. B. (1989). Organizational theories: Some criteria for evaluation. *The Academy of Management Review*, *14*(4), 469–515.

Barney, J. B., & Hesterly, W. (1996). Organizational economics: Understanding the relationship between organizations and economic analysis. In Glegg, S. R., Hardy, C., & Nord, W. R. (Eds.), *Handbook of organization studies* (pp. 115–147). London: Sage.

Barney, J. B. & Ouchi, W. G. (Eds.). (1986). *Organizational economics*. San Francisco, CA: Jossey-Bass.

Bergen, M., Dutta, S., & Walker, O. C. Jr. (1992). Agency relationships in marketing: A review of the implications and applications of agency and related theories. *Journal of Marketing*, *56*, 1–24.

Birnbaum, R. (1988). *How colleges work. The cybernetics of academic organization and leadership*. San Francisco, CA: Jossey-Bass.

Birnbaum, R. (2001). *Management fads in higher education. Where they come from, what they do, why they fail*. San Francisco, CA: Jossey-Bass.

Bleiklie, I. (1998). Justifying the evaluative state: New public management ideals in higher education. *European Journal of Education*, *33*(3), 299–316.

Bohte, J., & Meier, K. J. (2000). Goal displacement: Assessing the motivation for organizational cheating. *Public Administration Review*, *60*(2), 173–182.

Bowen, H. R. (1970). Financial needs of the campus. In Connery, R. H. (Ed.), *The corporation and the campus corporate support of higher education in the 1970s* (pp. 75–93). New York: The Academy of Political Science, Columbia University.

Bowen, H. R. (1980). *The costs of higher education. How much do colleges and universities spend per student and how much should they spend?* San Francisco, CA: Jossey-Bass Publishers.

Brewer, D. J., Gates, S. M., & Goldman, C. A. (2002). *In pursuit of prestige. Strategy and competition in U.S. higher education*. New Brunswick: Transaction Publishers.

Cave, M., Hanney, S., Henkel, M., & Kogan, M. (1997). *The use of performance indicators in higher education. The challenge of the quality movement* (3rd edn). London: Jessica Kingsley Publishers.

Clark, B. R. (1983). *The higher education system. Academic organization in cross-national perspective*. Berkeley, CA: University of California Press.

Clark, B. R. (1984). The organizational conception. In Clark, B. R. (Ed.), *Perspectives on higher education. Eight disciplinary and comparative views* (pp. 106–131). Berkeley, CA: University of California Press.

Dickoff, J., & James, P. (1992). A theory of theories: A position paper. In Nicoll, L. H. (Ed.), *Perspectives on nursing theory* (2nd edn) (pp. 99–111). Philadelphia, PA: J.B. Lippincott Co.

Donaldson, L. (1990). The ethereal hand: Organizational economics and management theory. *Academy of Management Review*, *1*(3), 369–381.

Donaldson, L. (1995). *American anti-management theories of organization. A critique of paradigm proliferation*. New York: Cambridge University Press.

Eggertsson, T. (1990). Economic behavior and institutions. Cambridge: Cambridge University Press.

Eisenhardt, K. (1989). Agency theory: An assessment and review. *Academy of Management Review*, *14*(1), 57–74.

Ferris, J. M. (1991). Contracting and higher education. *The Journal of Higher Education*, *62*(1), 1–24.

Garvin, D. A. (1980). *The economics of university behavior*. New York: Academic Press.

Geuna, A. (1999). *The economics of knowledge production. Funding and structure of university research*. Cheltenham, UK: Edward Elgar.

Gomez-Mejia, L. R., & Balkin, D. B. (1992). Determinants of faculty pay: An agency theory perspective. *Academy of Management Journal*, *35*(5), 921–955.

Gornitzka, Å., Stensaker, B., Smeby, J-C., & de Boer, H. (2004). Contract arrangements in the Nordic countries: Solving the efficiency/effectiveness dilemma? *Higher Education in Europe*, *29*(1), 87–101.

Hölttä, S. (1995). *Towards the self-regulative university*. Joensuu: University of Joensuu.

James, E. (1986). Cross-subsidization in higher education: Does it pervert private choice and public policy? In Levy, D. C. (Ed.), *Private education. Studies in choice and public policy* (pp. 237–257). New York: Oxford University Press.

James, E. (1990). Decision processes and priorities in higher education. In Hoenack, S. A. & Collins, E. L. (Eds.), *The economics of American universities. Management, operations, and fiscal environment* (pp. 77–106). Albany, NY: State University of New York Press.

Jensen, M. C., & Meckling, W. H. (1976). Theory of the firm: Managerial behavior, agency costs and ownership structure. *Journal of Financial Economics*, *3*(4), 305–360.

Johnes, J., & Taylor, J. (1990). *Performance indicators in higher education. UK universities*. Buckingham: SRHE & Open University Press.

Jongbloed, B., & Vossensteyn, H. (2001). Keeping up performances: An international survey of performance-based funding in higher education. *Journal of Higher Education Policy and Management*, *23*(2), 127–145.

Kells, H. R. (1994). Performance indicators for higher education: A critical review with policy recommendations. In Salmi, J. & Verspoor, A. M. (Eds.), *Revitalizing Higher Education* (pp. 174–208). Oxford: IAU Press/Pergamon.

Kezar, A. (2006). To use or not to use theory: Is that the question? In Smart, J. C. (Ed.), *Higher education: Handbook of theory and research,* Vol. XXI (pp. 283–344). Dordrecht: Springer.

Kiser, E. (1999). Comparing varieties of agency theory in economics, political science, and sociology: An illustration from state policy implementation. *Sociological Theory*, *17*(2), 146–170.

Kivistö, J. A. (2005). The government–higher education institution relationship: Theoretical considerations from the perspective of agency theory. *Tertiary Education and Management*, *11*(1), 1–17.

Kivistö, J. A. (2007). *Agency theory as a framework for the government–university relationship*. Tampere: Higher Education Group/Tampere University Press.

Laking, R. (2005). Agencies: Their benefits and risks. *OECD Journal on Budgeting*, *4*(4), 8–25.

Lane, J. E. (2005). *Agency problems with complex principals. State oversight of higher education: A theoretical review of agency problems with complex principals.* Paper presented at the Annual Conference of the Association for the Study of Higher Education (ASHE), Philadelphia, PA, November.

Lane, J. E., & Kivistö, J. A. (2008). The role of interests, information, and incentives in higher education: A review of principal–agent theory and its potential applications to study of higher education governance. In Smart, J. C. (Ed.), *Higher education: Handbook of theory and research,* Vol. XXIII (pp. 141–179). Dordrecht: Springer.

Levin, H. M., & Koski, W. S. (1998). Administrative approaches to educational productivity. *New Directions for Higher Education*, *103*, 9–21.

Liefner, I. (2003). Funding, resource allocation, and performance in higher education systems. *Higher Education*, *46*, 469–489.

Massy, W. F. (1996). Re-engineering resource allocation systems. In Massy, W. F. (Ed.), *Resource allocation in higher education* (pp. 15–47). Ann Arbor, MI: The University of Michigan Press.

McLendon, M. K. (2003). The politics of higher education: Toward an expanded research agenda. *Educational Policy*, *17*(1), 165–191.

Milgrom, P., & Roberts, J. (1992). *Economics, organization and management*. Upper Saddle River, NJ: Prentice Hall.

Mitnick, B. M. (1975). The theory of agency: The policing 'paradox' and regulatory behavior. *Public Choice*, *24*, 27–42.

Moe, T. M. (1984). The new economics of organization. *American Journal of Political Science*, *28*(4), 739–777.

Ortmann, A., & Squire, R. (2000). A game–theoretic explanation of the administrative lattice in institutions of higher learning. *Journal of Economic Behavior & Organization*, *43*, 377–391.

Perrow, C. (1986). *Complex organizations. A critical essay* (3rd edn). New York: Random House.

Rose-Ackerman, S. (1978). *Corruption. A study in political economy*. New York: Academic Press.

Ross, S. A. (1973). The economic theory of agency: The principal's problem. *American Economic Review*, *63*(2), 134–139.

Schmidtlein, F. A. (2004). Assumptions commonly underlying government quality assessment practices. *Tertiary Education and Management*, *10*(4), 263–285.

Tosi, H. L. Jr., & Gomez-Mejia, L. R. (1989). The decoupling of CEO pay and performance: An agency theory perspective. *Administrative Science Quarterly*, *34*(2), 169–189.

Trow, M. A. (1996). Trust, markets and accountability in higher education: A comparative perspective. *Higher Education Policy*, *9*(4), 309–324.

Vedder, R. (2004). *Going broke by degree. Why college costs too much?* Washington, DC: The AEI Press.

Vining, A., & Globerman, S. (1999). A conceptual framework for understanding the outsourcing decision. *European Management Journal*, *17*(6), 645–654.

van Vught, F. A. (1997). The effects of alternative governance structures. In Steunenberg, B. & van Vught, F. A. (Eds.), *Political institutions and public policy. Perspectives on European decision making* (pp. 115–137). Dordrecht: Kluwer Academic Publishers.

Waterman, R. W., & Meier, K. J. (1998). Principal–agent models: An expansion. *Journal of Public Administration Research and Theory*, *8*(2), 173–202.

Weick, K. E. (1989). Theory construction as disciplined imagination. *Academy of Management Review*, *14*(4), 516–531.

Whicker, M. L. (1997). An economic perspective of academic tenure. *PS: Political Science and Politics*, *30*(1), 21–25.

Whynes, D. K. (1993). Can performance monitoring solve the principal–agent problem? *Scottish Journal of Political Economy*, *40*(4), 434–446.

Williams, G. (1995). The 'marketization' of higher education: Reforms and potential reforms in higher education finance. In Dill, D. D. & Sporn, B. (Eds.), *Emerging patterns of social demand and university reform: Through a glass darkly* (pp. 170–193). Oxford: IAU Press/ Pergamon.

## [Alexander Pepper](#) and Julie Gore

# Behavioral agency theory: new foundations for theorizing about executive compensation

## Article (Accepted version)
## (Refereed)

# Behavioral Agency Theory: New Foundations for Theorizing about Executive Compensation

Alexander Pepper

*The London School of Economics and Political Science*

Julie Gore

*University of Surrey*

*Corresponding author: Alexander Pepper, Department of Management, The London School of Economics and Political Science, Houghton Street, London, WC2A 2AE, United Kingdom*

*Email: a.a.pepper@lse.ac.uk*

# ABSTRACT

This paper describes new micro-foundations for theorizing about executive compensation, drawing on the behavioral economics literature and based on a more realistic set of behavioral assumptions than those which have typically been made by agency theorists. We call these micro-foundations "behavioral agency theory". In contrast to the standard agency framework, which focuses on monitoring costs and incentive alignment, behavioral agency theory places agent performance at the center of the agency model, arguing that the interests of shareholders and their agents are most likely to be aligned if executives are motivated to perform to the best of their abilities. We develop a line of argument first advanced by Wiseman and Gomez-Mejia (1998), and put the case for a more general reassessment of the behavioral assumptions underpinning agency theory. A model of economic man predicated on bounded rationality is proposed, adopting Wiseman and Gomez-Mejia's assumptions about risk preferences, but incorporating new assumptions about time discounting, inequity aversion and the trade-off between intrinsic and extrinsic motivation. We argue that behavioral agency theory provides a better framework for theorizing about executive compensation, an enhanced theory of agent behavior and an improved platform for making recommendations about the design of executive compensation plans.

**Keywords:** agency theory; behavioral theory; compensation, bonuses and benefits; motivation; top management teams

# INTRODUCTION

Agency theory has been a major component of the economic theory of the firm since the publication of formative work by Spence and Zeckhauser (1971), Alchian and Demsetz (1972), Ross (1973) and Jensen and Meckling (1976). It has also become the dominant theoretical framework for academic research on executive compensation (Bratton, 2005). The literature on senior executive reward is now very extensive, drawing on a variety of scholarly traditions, including economics, law, organization studies, accounting and finance. In addition to the agency approach, theoretical frameworks include tournament theory (Lazear & Rosen, 1981), human capital theory (Combs & Skills, 2003), the managerial-power hypothesis (Bebchuk, Fried & Walker, 2002), institutional theory (Balkin, 2008), political theories (e.g., Ungson & Steers, 1984) and theories about fairness (e.g., Wade, O'Reilly & Pollock, 2006). Literature reviews and summaries are provide by Gomez, Meija & Wiseman (1997), Devers, Cannella, Reilly and Yoder (2007) and Gomez-Mejia, Berrone and Franco-Santos (2010: 117-140). Denvers et al (2007) note that behavioral research is a relative new feature of the literature on senior executive reward.

That agency theory has shortcomings has been apparent for some time. Most notably, given Jensen's role as a leading agency theorist, empirical work carried out by Jensen and Murphy (1990) failed to establish a conclusive link between CEO pay and stock price performance[1] Ten years later, in a meta-analysis of 137 empirical studies, Tosi, Werner, Katz and Gomez-Mejia (2000) similarly found that incentive alignment as an explanatory agency construct for CEO pay was at best weakly supported by the evidence. More recently, Frydman and Jenter (2010) have argued, based on a review of US executive compensation data covering the period 1936 to 2005, that neither optimal contracting (agency theory) nor the managerial power hypothesis is fully consistent with the available evidence. Roberts, another agency theorist, has commented that agency theory performed poorly during the

financial crisis and has reported various situations where strong incentives are evidently not optimal, as agency theory implies (Roberts, 2010). These include when good measures of an agent's effort or performance are not available, when multi-tasking is required, and when cooperation between different agents is necessary, all common situations where top management teams are concerned. Roberts puts forward arguments in favor of implementing *weak* rather than *strong* incentives in such circumstances.

This paper proposes a new version of agency theory which provides a better explanation of the connection between executive compensation, agent performance, firm performance and the interests of shareholders. We call this "behavioral agency theory", developing a line of argument first advanced by Wiseman and Gomez-Mejia (1998), who proposed that the normal risk assumptions of agency theory should be varied to incorporate ideas from prospect theory (Kahneman & Tversky, 1979; Tversky & Kahneman, 1992). Sanders and Carpenter (2003) have subsequently adopted a behavioral agency perspective in their examination of stock repurchase programs and a summary of the literature using the behavioral agency framework is provided by Finkelstein, Hambrick and Cannella (2009). Rebitzer and Taylor (2011) provide a general examination of behavioral approaches to agency and labor markets in the 4[th] edition of Ashenfelter and Card's influential handbook on labor economics. However, a settled theory and agreed terminology for the behavioral agency model does not yet exist. In contrast to the standard agency framework, which focuses on monitoring costs and incentive alignment, behavioral agency theory places agent performance and work motivation at the center of the agency model, arguing that the interests of shareholders and their agents are most likely to be aligned if executives are motivated to perform to the best of their abilities, given the available opportunities. Behavioral agency theory builds on four constructs which have been identified as key factors affecting behavior by behavioral economists (Camerer, Loewenstein & Rabin, 2004). These are: (1) loss aversion and

3

reference dependence; (2) preferences relating to risky and uncertain outcomes; (3) temporal discounting; and (4) fairness and inequity aversion. It incorporate a theory (crowding out) relating to the trade-off between intrinsic and extrinsic motivation (Frey & Jegen, 2001; Sliwka, 2007). It also introduces goal-setting theory (Locke & Latham, 1984, 1990) to the agency model, on the basis that it represents a pragmatic way of contracting between principal and agent.

The paper proceeds as follows: it begins by describing agency theory's main elements and underlying assumptions, before reviewing the limitations of positive agency theory as an explanation of the relationship between senior executives and shareholders, and reconceptualizing what is meant by economic man (i.e., homo economicus of neoclassical economics). It continues with an explanation of the behavioral agency model, describing the main component systems and commenting in some detail on the significance of motivation, risk, time discounting, inequity aversion and goal setting. It examines the relationship between job performance and firm performance, discusses ways in which behavioral agency theory departs from standard agency theory, and considers the implications of behavioral agency theory for compensation design, before concluding.

## POSITIVE AGENCY THEORY

Positive agency theory[2], the standard model of agency which we consider in this paper, has been extensively used as a basis for theoretical and empirical work by management scholars and organization theorists (e.g., Eisenhardt, 1989), as well as being widely applied in examining research questions relating to executive compensation (e.g., Tosi & Gomez-Mejia, 1989). It argues that the firm is a special case of the theory of agency, that a firm provides a nexus for a complex set of contracts, both written and unwritten, between various parties, and that agency costs are generated as a result of the different interests and contractual

arrangements between owners and top managers (Alchian & Demsetz, 1972; Jensen, 1983; Jensen & Meckling, 1976). The underlying assumptions are that organizations are profit seeking, that agents are both rational and rent seeking, and that there is no non-pecuniary agent motivation[3]. It is further assumed that principals are risk neutral, because they can balance their portfolios, that agents are risk averse, because the potential wealth effects of the employment relationship are significant, that an agent's utility is positively contingent on pecuniary incentives and negatively contingent on effort, and that time preferences are calculated mathematically according to an exponential discount function (Jensen, 1998). It is postulated that effort and motivation increase monotonically with additional reward[4]. The pay-effort function is visualized as a straight line with a positive gradient proceeding from bottom left to top right.

Efficiency is the main criterion for assessing the success or otherwise of programs under agency theory. Agency theory focuses on the costs of the potential conflict of interest between principals and agents, referred to as "agency costs". Jensen and Meckling define agency costs as the sum of the monitoring expenditures of the principal, the bonding expenditures of the agent, and the residual loss in welfare experienced by the principal as a result of the divergence of interests between the principal and the agent (Jensen & Meckling, 1976). Jensen subsequently offers a broader definition, describing agency costs as "the sum of the costs of structuring, bonding and monitoring contracts between agents…[which]…also include the costs stemming from the fact that it does not pay to enforce all contracts perfectly" (Jensen, 1983: 331)[5]. Agency costs are thus a special case of transaction costs (in a Coasian sense) in their internal (intra-firm) rather than external (intra-market) form. Positive agency theory proposes that principals can mitigate agency costs by establishing appropriate incentive contracts and by incurring monitoring costs. This is formalized by Eisenhardt in two propositions - first, in respective of incentives: "when the contract between

the principal and agent is outcome based, the agent is more likely to behave in the interests of the principal" (Eisenhardt, 1989: 60); secondly, in respect of monitoring: "when the principal has information to verify agent behavior, the agent is more likely to behave in the interests of the principals" (Eisenhardt, 1989: 60).

## BEHAVIORAL AGENCY THEORY

Behavioral agency theory argues that the model of economic man which forms the micro-foundations of agency theory is too simplistic. It proposes a reconceptualization, developing a new model which assumes bounded rationality[6], recognizes the importance of agents' human capital (taking this to be a function of ability and work motivation) and allows for departures from the rational choice model when it comes to loss, risk and uncertainty aversion, time discounting, inequity aversion and the trade-off between intrinsic and extrinsic motivation. It proposes that the standard agency theory model of man should be modified in a number of ways. The first modification relates to agent performance and work motivation. Agency theory places less emphasis on the objective of motivating agents to perform to the best of their ability than it does on aligning the interests of agents and principals. Leibenstein (1966) argues that, given the importance of what is now called human capital, motivation (in particular, intrinsic motivation) cannot be ignored in the economic calculus. Pratt and Zeckhauser (1985) make the same case for agency theory. Behavioral agency theory argues that maximizing agent performance should be a key objective of the principal-agent relationship and that the importance of the agent's work motivation, including intrinsic motivation, should not be underestimated. It challenges the idea that intrinsic and extrinsic motivation are either independent or additive, arguing instead that contingent monetary rewards might actually cause a reduction in intrinsic motivation (see Deci & Ryan, 1985).

Frey and Jegen (2001), following a line of scholarly thinking that dates back to Lepper and Greene (1978), have described this phenomenon as "crowding-out"(see also Sliwka, 2007) .

The second modification relates to risk and uncertainty[7].  Behavioral agency theory assumes that senior executives are primarily loss averse and only secondarily risk averse (Wiseman & Gomez-Mejia, 1998).  Gains and losses are calculated by each individual agent in relation to a reference point which he or she subjectively determines.  Risk preferences differ in gains and losses, resulting in an "S-shaped" value function, with losses looming larger than gains.  This means that, below a reference point, agents will be loss averse, resulting in an increase in his or her appetite to take short term risk.  Above the reference point agents will generally be risk averse, but decision weights will vary depending on subjective probability assessment; for example, small probabilities are over-weighted and large probabilities are under-weighted.

The third modification relates to time preferences.  In behavioral agency theory it is assumed that agents discount time according to a hyperbolic discount function, rather than exponentially, as is the case with financial discounting (Ainslie, 1991; Ainslie & Haslam, 1992).  This means that future rewards are heavily discounted and allows for the possibility of preference reversals. Actual average discount rates vary between individuals and must be determined empirically.

The fourth modification relates to an agent's perceptions of equitable compensation.  If agents feel that their inputs, the effort and skills which they put in to their work, are fairly and adequately rewarded by outputs, the tangible and intangible rewards from employment, then the agents will be happy in their work and motivated to continue to contribute at the same or at a higher level (Adams, 1965).  However, if the relationship between inputs and outputs is not proportionate, then an agent will become dissatisfied and hence demotivated.  In this model the agent's equity benchmark is subjectively determined according to market norms

and personal referents. Fehr and Schmidt (1999: 819) call this phenomenon "inequity aversion". As is the case with risk and time discounting, we anticipate that actual levels of inequity aversion will vary between individuals and must be determined empirically.

Table 1 summarizes the assumptions about the characteristics of economic man which provide the foundations of agency theory and compares them with the way in which behavioral agency theory reconceptualizes the model. An important early conclusion which can be drawn is that an agent's perception of the (subjectively-calculated) value of an incentive award will typically be less than the award's (objectively-calculated) economic value. This clearly has implications for the way that incentive contracts are designed.

```
┌─────────────────────────────────────────────┐
│                                             │
│            TABLE 1 ABOUT HERE               │
│                                             │
└─────────────────────────────────────────────┘
```

## Assessment Criteria and Unit of Analysis

Behavioral agency theory proposes that it is necessary to use both effectiveness and efficiency as yardsticks for judging agent activity. By adopting effectiveness as well as efficiency as criteria for assessment we follow a long line of management theorists dating back to Barnard (1938 |1968)[8]. Simon (1945 |1997) pointed out that the terms "effectiveness" and "efficiency" were considered to be almost synonymous until the end of the 19[th] century and were generally thought to mean the power to accomplish the purpose intended; however, the meanings of the two words subsequently diverged. Efficiency came to be defined, firstly in engineering and subsequently in economics, business, and management, in terms of the relationship between inputs and outputs. In this paper we use the terms efficiency and effectiveness in the following way: on the one hand, an action, event, plan, policy or program is considered to be *efficient* if it causes inputs to be minimized for a

8

given set of outputs or outputs maximized for a given set of inputs; on the other hand, an action, event, plan, policy or program is considered to be *effective* if it is capable of achieving its intended objectives, what the objectives are being exogenous to the theory. We contend that it is necessary for management scholars to adopt both criteria in order to provide a complete and accurate evaluation of management policies, plans and programs. Taking executive compensation as an example, a compensation plan might be effective and efficient (i.e., achieve its objectives of motivating top managers and aligning the interests of managers and shareholders, doing so in such a way that costs are minimized), effective but not efficient (i.e., achieve its objectives but in a way that is more costly than necessary), or neither effective nor efficient (i.e., fail to achieve its objectives at the same time as being costly). However, we argue that it makes no sense to describe a management plan or program as efficient but not effective. The concept of effectiveness is already implied by the concept of efficiency; a lower cost (or indeed no cost at all) could otherwise be incurred while still failing to achieve the desired objectives.

An important premise of behavioral agency theory, consistent with the top management team or "upper echelons" approach (Hambrick & Mason, 1984) is that senior executive teams have a major impact on firm performance. We define "top management team" (and hence "top manager") as the group of very senior executives who are responsible for defining and executing a firm's strategy, who through their actions are capable of affecting the company's profits, share price, reputation and market positioning (Carpenter, Geletkanycz & Sanders, 2004). This group, which includes the chief executive officer (CEO), the chief operating officer (COO), the chief financial officer (CFO), divisional heads and other heads of function, is sometimes referred to as the "management board", "operating board", "executive committee" or "general management committee". Changing trends in corporate governance mean that, while historically these individuals would have been executive directors, it is

<div align="center">9</div>

increasingly common in many countries to find only the CEO and CFO on the main board, while all the key senior executives sit on the executive committee, or equivalent (Pepper, 2006). By defining top managers in this way, this part of behavioral agency theory becomes, in a sense, tautological (corporate performance is, in part, a function of the performance of top managers; top managers are those individual agents who are able to influence corporate performance). However, this is the type of "useful tautology" which Jensen (1983: 330-331) points out is a necessary part of the process of theory development; nor does its inclusion in behavioral agency theory mean that this part of the theory becomes in practice irrefutable - it might be demonstrated in certain cases that top managers are not in practice able to have a significant impact on firm performance.

Unlike upper echelons theory, which takes the top management team as the primary unit of analysis (Hambrick & Mason, 1984) behavioral agency theory focuses on the behaviors, interests and actions of individual top managers or agents. Following Boxall and Purcell (2003) we model an agent's performance as a manager of a large firm as a function of his or her ability, motivation and opportunity. Agents will perform if they have the ability (the necessary knowledge, skill and aptitude), the motivation (intrinsic and extrinsic), and the right opportunities (including the necessary work structures and business environment); formally:

$$P_{\bar{a}} = f(A, M, O) \tag{1}$$

where $P_{\bar{a}}$ stands for the job performance of the agent, A stands for ability, M stands for motivation or "motivational force", after Lewin (1938), and O stands for the agent's opportunity set.

Boxall and Purcell conceptualize ability in much the same way that Becker (1993) conceptualizes human capital, i.e., in terms of knowledge, skills, health, value and habits. Leibenstein (1966) comments on the importance of motivation to human capital. The

significance for behavioral agency theory is that a competent agent must be properly motivated in order to ensure optimal performance (Pratt & Zeckhauser, 1985), meaning in this context the point where efficiency is maximized subject to any effectiveness constraints, and effectiveness is maximized subject to any efficiency constraints. Thus we define human capital in this article as "motivated ability" rather than merely as a function of education and experience.

In this paper we focus on the role of motivation in influencing the job performance of agents. For the purposes of the current paper, we take ability, which has its roots in the learning and development and human capital literatures, and opportunity, which can be traced to the leadership and strategy literatures, as given.

## Motivation

The theory of work motivation most commonly used in investigations into the motivational impact of pecuniary incentives is expectancy theory (Vroom, 1964). According to expectancy theory, motivational force is a function of expectancy (the strength of belief or subjective probability that an action i will lead to a particular outcome j), instrumentality (the degree to which a first outcome j will lead to a second outcome k), and valence (the preference which an individual has for the second outcome k)[9]. Expectancy theory is essentially concerned with extrinsic, rather than intrinsic or total motivation. Thus expectancy theory can be formally represented as:

$$X_i = f(E_{ij}, I_{jk}, V_k) \tag{2}$$

where $X_i$ is the extrinsic motivational force to perform act i, $E_{ij}$ is the strength of expectancy that act i will be followed by outcome j, $I_{jk}$ is the instrumentality of outcome j for attaining outcome k, and $V_k$ is the valence of outcome k. Expectancy theory thus describes a cognitive process and is distinct from many of the other standard theories of motivation, especially

11

theories based on needs, drives and learned behaviors, which seek to explain the

psychological content of motivation

Steel and König (2006) have proposed a modified version of expectancy theory which

they call "temporal motivation theory". It postulates that motivation can be understood in

terms of valence and expectancy[10], weakened by delay, influenced by risk and uncertainty,

with different valences for gains and perceived losses. Temporal motivation theory brings

expectancy theory together with prospect theory (Kahneman & Tversky, 1979; Tversky &

Kahneman, 1992) and hyperbolic discounting (Ainslie, 1991; Ainslie & Haslam, 1992).

Reducing Steel and König's formula down to its minimal form gives:

$$X_i = \left\{ \frac{E_{ik}^{pt} \times V_k^{pt}}{1 + \delta t} \right\} \tag{3}$$

where $X_i$ is again the extrinsic motivational force to perform act i, $E_{ik}^{pt}$ is the expectancy that

act i will lead, via j, to outcome k, $V_k^{pt}$ is the valence for outcome k, $\delta$ is the personal

discount factor for the delay between act i and outcome k, and t represents the time-lag.

Expectancy and valence are both calculated in accordance with prospect theory. The main

implications of this are that probabilities and decision weights are determined subjectively

and valence is affected by risk perception: in particular, valences will differ significantly

depending on whether gains or losses are expected[11]. Time effects are determined by a

hyperbolic discount function after Ainslie (1991) rather than the more conventional

exponential discounting function used in financial theory. This means that, in Steel and

König's revised expectancy model, the valence which an agent attaches to k takes into

account risk and uncertainty, as well as being discounted for any time delay between the

occurrence of act i and outcome k.

Positive agency theory places less emphasis on the objective of motivating agents than it

does on alignment of the interests of agents and principals. Deci and Ryan (1985) point out

that there are two distinct forms of motivation, intrinsic motivation, where an agent performs

an activity for its inherent satisfaction rather than because of some separable consequence, and extrinsic motivation, where an activity is carried out because of its instrumental value. Kreps (1997) argues that it is not necessary to postulate the concept of intrinsic motivation on the basis that what is called intrinsic motivation may in fact be no more than a series of vaguely defined extrinsic motivators. Besley and Ghatak (2005) contend, on the contrary, that there is such a thing as a motivated agent whose economic behavior is affected by intrinsic motivation, but their argument is directed towards employees of public sector or non-profit organizations whose activities coalesce around a "mission". Deci and Ryan (1985) argue that the importance of intrinsic motivation should not be underestimated. They challenge the idea that intrinsic and extrinsic motivation are either independent or additive, arguing instead that contingent monetary awards might actually cause a reduction in intrinsic motivation. Boivie, Lange, McDonald and Westphal (2011) point out how, in the case of CEOs, high organizational identification, which may be associated with intrinsic motivation, can help to reduce agency costs. Frey and Jegen (2001) and Sliwka (2007) postulate that in some cases extrinsic rewards can "crowd-out" intrinsic motivation, particularly if monetary incentives are badly designed. They argue for a strong form of crowding-out whereby an increase in extrinsic reward leads to an overall reduction in total motivation. A weaker form of crowding-out, whereby the level of total motivation is maintained only if the increase in extrinsic reward more than compensates for the reduction in intrinsic motivation, can alternatively be postulated.

Following Deci and Ryan, (1985), the relationship between intrinsic and extrinsic motivation can be stated formally as follows:

$$M_i = f\,(N_i,\, X_i) \tag{4}$$

where $M_i$ is an agent's total motivational force, $N_i$ is the agent's intrinsic motivation, and $X_i$ is the agent's extrinsic motivation. $M_i$, $N_i$, and $X_i$ can be thought of in terms of stimuli, actions

or outcomes i.e., $M_i$ represents motivation resulting from i, where i is a stimulus or bundle of stimuli, an action or package of actions, an outcome or collection of outcomes. However, the relationship between $N_i$ and $X_i$ is neither linear nor additive (Deci & Ryan, 1985). In a dynamic sense, when changes in incentives occur, there is evidently a trade-off of some kind between the two types of motivation. Whether this is more accurately described by the strong crowding-out conjecture, where a change in extrinsic motivation as earnings increase from e to g, $(+\Delta X_{eg})$ leads to a decrease in intrinsic motivation $(-\Delta N_{eg})$ such that $\Delta N_{eg} > \Delta X_{eg}$ and $M_e > M_g$, or by the weak crowding-out conjecture, such that $\Delta N_{eg} = \Delta X_{eg}$ and $M_e = M_g$, can only be determined empirically. This argument leads to the first two research propositions:

> *Proposition 1a: (The weak crowding-out conjecture) Above a certain level of compensation (represented by inflection point $\lambda_1$ on the agent's pay-effort curve) intrinsic motivation will decrease as compensation increases, such that the rate of increase of total motivation will diminish and will eventually, at a higher level of compensation (represented by inflection point $\beta$ on the agent's pay-effort curve), reach zero.*

> *Proposition 1b: (The strong crowding-out conjecture) If compensation continues to increase above the higher level of compensation represented by inflection point $\beta$ on the agent's pay-effort curve, then total motivation will start to decline as intrinsic motivation is crowded out by extrinsic rewards.*

Compensation comprises the sum of all incentives and rewards, pecuniary and non-pecuniary, arising from the agency relationship. The difference between incentives and rewards is that incentives are determined ex ante (i.e., prior to performance, thus encouraging agents to act in a particular way) whereas rewards are determined ex post.

**Risk**

A standard assumption of agency theory is that agents are risk averse. According to behavioral agency theory this is an oversimplification. We argue that extrinsic motivation and agent behavior are significantly affected by the agent's risk profile and that a more sophisticated model of risk and uncertainty is accordingly required. Behavioral agency theory postulates, after Wiseman and Gomez-Meija (1998), who in turn cite Kahneman and Tversky (1979) and Tversky and Kahneman (1992), that agents are primarily loss averse and consequently, contrary to one of the standard assumptions of agency theory, may actually have a high propensity to take short term risks below a certain level of compensation representing the point where perceived gains become perceived losses. Above this gain/loss inflection point, agents will generally be risk averse, but small probabilities are typically over-weighted and large probabilities are typically under-weighted. The gain/loss inflection point is itself context dependent and a matter of individual differences. In particular it is contingent upon the agent's perception of his or her individual compensation endowment which comprises their actual current compensation, enhanced to the extent of future incentives which are expected to be received with a reasonable degree of certainty. For example, a future bonus which is guaranteed or otherwise strongly anticipated based, say, on the pattern of past bonus payments, would be taken into account in the current compensation endowment, albeit discounted for future payment. In a similar way, an agent with underwater options (where the current stock price is below the option strike price) may regard this as representing a loss on his or her current compensation endowment.

This enables us to advance, following Wiseman and Gomez-Meija (1998), two further propositions:

> *Proposition 2a: Below a level of compensation (represented by inflection point $\lambda_2$ on the agent's pay-effort curve) agents are loss averse.*

15

*Proposition 2b: Above a level of compensation (represented by inflection point λ₂ on the agent's pay-effort curve) agents are risk averse.*

## Time Discounting

Positive agency theory assumes that time differences can be accounted for by the type of conventional exponential discount function used in finance theory. However, behavioral economists have identified a series of anomalies in the way that individuals account for time, including preference reversal and weakness of will (undertaking actions which in the short term are pleasurable, but which agents know to be detrimental to their well-being in the long term). Ainslie (1991) explains these anomalies by arguing, based on experimental evidence, that his subjects discount future events hyperbolically so that the implied discount rate varies over time, rather than exponentially, which would require a constant discount rate. That economic agents typically discount time hyperbolically is generally accepted as the norm by behavioral economists (Frederick, Loewenstein & O'Donoghue, 2002; Graves & Ringuest, 2012). Steel and König (2006) argue that expectancy theory must take into account time differences as compensation (outcome k in Equation 3 above) may not be received until sometime after the action which leads to the payment (act i). They also argue that time differences should be accounted for using a hyperbolic discount function. Accordingly, we postulate that an agent's extrinsic motivation is affected by time discounting, calculated on a hyperbolic discount basis, as set out in the next proposition:

*Proposition 3: Agents discount future compensation according to a hyperbolic discount factor such that the average discount rate δ is significantly greater than the equivalent financial discount rate.*

**Inequity Aversion**

Behavioral agency theory postulates that motivational force is affected by inequity aversion, based on equity theory (Adams, 1965). It is widely recognized that an individual's satisfaction with his or her compensation depends not just upon buying-power, but also on how their compensation compares with the total rewards of salient others (Shafir, Diamond & Tversky, 1997). Akerlof (1982) postulates the fair-wage hypothesis according to which workers have a conception of a "fair-wage" such that, if actual earnings are less than the fair-wage, then only a corresponding fraction of normal effort will be supplied. According to Adams (1965) people seek a fair balance between what they put into their jobs (including energy, commitment, intelligence and skill – collectively "inputs") and what they get out (including financial rewards, recognition, and opportunities for personal growth – collectively "outputs"). Agents form perceptions of what constitutes an appropriate balance between inputs and outputs by comparing their own situations with those of other people in accordance with the ratio $\{O_{\bar{a}}/I_{\bar{a}}\} : \{O_r/I_r\}$ (which we refer to below as the "Adams' ratio") where $O_{\bar{a}}$ is the agent's outputs e.g., their compensation, $I_{\bar{a}}$ is the agent's inputs e.g., their skills and effort, $O_r$ is the outputs of the agent's referents and $I_r$ is the referents' inputs. Referents may be internal (peers, immediate subordinates, immediate superiors) or external (people doing equivalent jobs in other organizations). If agents feel that their inputs are fairly and adequately rewarded by outputs, the equity benchmark being subjectively perceived from market norms and other reference points, then they will be happy in their work and motivated to keep contributing at the same or a higher level. However, if the relationship between inputs and outputs is not proportionate, such that $\{O_{\bar{a}}/I_{\bar{a}}\} < \{O_r/I_r\}$, then the agent will become dissatisfied and hence demotivated. "Inequity aversion", as Fehr and Schmidt (1999: 819) call this phenomenon, is translated by Michelman into economic terms as "demoralization costs" (Michelman, 1967: 1214). Gomez-Mejia and Wiseman (1997) argue

that inequity aversion applies equally to senior executives as to other workers. Accordingly, we generate the next proposition:

*Proposition 4: Individual agents will determine a level of compensation (represented by inflection point λ₃ on the agent's pay-effort curve) by reference to the compensation of a class of significant referents, such that the agent will tend to be dissatisfied and hence demotivated if his or her actual earnings are less than λ₃.*

It is important to note for Proposition 4 that an individual agent's assessment of relative compensation levels will take account of his or her perception of their contribution in comparison with that of his or her referents, in accordance with what we have described above as the "Adams' ratio", according to which individuals seek to balance perceived relative inputs and outputs (Adams, 1965).

## Goal Setting, Contracting and Monitoring

We turn now to goal setting, contracting and monitoring. We argue that these activities should be seen as integral to behavioral agency theory: goal setting and monitoring are important factors in legal contracting, which is a key element in the relationship between principal and agent (Grossman & Hart, 1983; Hart, 1995); they have also been demonstrated to be an important component of agent motivation (Locke & Latham, 1984, 1990). Goal setting theory postulates a strong connection between goals, commitment and performance. Goals must be specific, difficult, attainable, and self-set or explicitly agreed to for the motivational affect to be maximized. Much of the empirical work supporting goal setting theory has been carried out in an industrial context (e.g., with loggers, truck drivers and word processing operators). Nevertheless, behavioral agency theory postulates that many of the features of goal setting theory are generalizable to senior executives. Locke and Latham (2002) make three points which are particularly pertinent to agency relationships. First, they

argue that monetary incentives enhance goal commitment but have no substantive effect on motivation unless linked to goal setting and achievement. Secondly, they explain, through a model which they call the "high performance cycle", how goal setting and achievement together lead to high performance, in turn leading to rewards, high job-satisfaction and self-efficacy. Thirdly, they suggest a possible connection with prospect theory, both theories stressing the importance of reference points in cognition.

One of the main problems with the relationship between principals and agents which has been identified by agency theorists is that agency contracts are inevitably incomplete (Grossman & Hart, 1983; Hart, 1995). If principals were able to specify completely all that they required of their agents, then there might be no need for incentive contracts to align the interests of principals and agents - monitoring of actions and outcomes might suffice. However, in practice there are limits on knowledge and cognition. One of the reasons that principals employ agents is for the agents' expertise. An agent who is more knowledgeable about the matters which are to be specified in a contract may be able to second-guess the principal during and after contract negotiation to the agent's advantage and the principal's detriment. There are also dynamic constraints. Over the course of time the business environment which provides the backdrop for the agency contract inevitably changes. Actions which are contractually required of the agent when a contract is negotiated may cease to be appropriate at a later date because of environmental changes, and other actions which could not have been anticipated ex ante may subsequently become necessary ex post. It is contractual uncertainties of this kind that Roberts (2010) is referring to when he advocates the merits of weak rather than strong incentives in agency relationships. Goal setting, especially when it involves discussions between principal and agent about the appropriate level of objectives, is a pragmatic way of contracting, given limits on knowledge and cognition. It is also a signaling mechanism, indicating to one of the parties in an exchange relationship, the

agent, what is required by another party, the principal. Spence (1973) has shown how signaling mechanisms of this kind form an important part of economic exchange in the context of employment. Thus, goal setting, monitoring and reward, as part of a regular high performance management cycle, provide a way of improving the quality of contracting in a manner which helps to enhance rather than undermine agent motivation. This leads to a further proposition:

> *Proposition 5a: The existence in a firm of a system of goal setting, monitoring, and linked rewards and incentives for agents who are members of the top management team is positively correlated with agent performance and work motivation.*

Some care is required, however. First, it would not possible to specify in a performance contract a full set of the objectives which would be necessary to cover all possible situations that might arise during the course of a performance cycle. According to the principle of requisite variety, a control system requires a response mechanism for every exogenous shock which it might face (Ross Ashby, 1956 |1976). Top managers face great complexity in their work and it would not be possible to anticipate every possible exogenous shock in a performance contract, nor to specify fully all the requirements of the job (Mintzberg, 1997, 2009). Arrow (1985) notes how unrealistic such a complex fee function would be. Secondly, the knowledge constraints of the bounded rationality assumption place cognitive limits on an agent's ability to assimilate and understand complex goals and performance criteria. This in turn leads to Proposition 5b, which is consistent with the conclusions reached by Roberts (2010), described in the introduction:

> *Proposition 5b: Weak incentives are a more effective and efficient way of motivating agents than strong incentives.*

**Agent's Job Performance and Work Motivation Cycle**

The various elements of the subsystem which models agent job performance and work motivation are summarized in Figure 1. This figure illustrates the trade-off between intrinsic and extrinsic motivation (the subject of Propositions 1a and 1b), the roles played by risk (Propositions 2a and 2b), time discounting (Proposition 3) and inequity aversion (Proposition 4). The goal setting, contracting and monitoring process (Propositions 5a and 5b) are illustrated, along with the integral feedback mechanism. Two further propositions, developed later in the paper, are also represented.

INSERT FIGURE 1 ABOUT HERE

In our analysis of the trade-off between intrinsic and extrinsic motivation, risk, and inequity aversion we have identified three compensation inflection points on the agent's pay-effort curve: $\lambda_1$ which is critical to the trade-off between intrinsic and extrinsic motivation; $\lambda_2$ which determines where an individual's risk appetite changes from loss aversion to risk aversion; and $\lambda_3$ which acts as the reference point for comparisons with salient others in the context of inequity aversion. As we have explained, $\lambda_1$, $\lambda_2$ and $\lambda_3$ are critical points in the various sub-systems. In our representation of an agent's pay-effort function in Figure 2 we make the assumption that these three inflection points are identical for any one individual agent. There is support for this assumption in the argument advanced by Deci and Ryan (1985) that the psychological sub-systems for intrinsic motivation, risk, and inequity aversion converge upon a common psychological state in which cognitive, affective and conative

21

variables are optimally aligned.  However, we assume the equality of the three inflection

points largely for mathematical convenience.  In practice, even if there is a linear range

between an upper inflection point ($\lambda_1$) and a lower inflection point ($\lambda_2$ and $\lambda_3$), or a plane with

three separate inflection points, the main argument, which is that there is a set of values for

which an agent's pay-effort ratio is optimized, would not be undermined - the range of

possible outcomes would simply be expanded.

Intrinsic motivation is represented in Figure 2 by the $\varepsilon = f(N_i)$ curve and extrinsic

motivation by the $\varepsilon = f(X_i)$ curve.  By superimposing the extrinsic motivation curve on top of

the intrinsic motivation curve, we generate the total motivation or $\varepsilon = f(N_i, X_i)$ curve.   This

runs parallel to the extrinsic motivation curve until total compensation reaches $\omega^*$, at which

point crowding out sets in, intrinsic motivation starts to decline and the rate of increase of

total motivation slows accordingly.

FIGURE 2 ABOUT HERE

By assuming the equality of $\lambda_1$, $\lambda_2$ and $\lambda_3$ (i.e., $\lambda$) we infer that there is a preferred level of

pay at which point the relationship between agent motivation ($\varepsilon^*$) and total compensation

($\omega^*$) is optimized, subject to constraints for risk, time discounting and inequity aversion.

This is the point when an agent's effort-to-pay ratio is maximized, such that the gradient of

the total motivation curve is equal to 1.   It implies that there is a set of first best

compensation strategies, being combinations of fixed and variable pay, contingent and

discretionary bonuses, and short-term and long term incentives: formally, that $\lambda$ is represented

by the set $\{\sigma_1, \sigma_2, \ldots \sigma_n\}$ where $\sigma$ represents a compensation mix with a unique combination of

22

464

fixed, variable, contingent, discretionary, current and deferred rewards.  If in practice there was a linear range between an upper and lower inflection point, or a plane with three separate inflection points, then this would simply increase the set of first best pay combinations. Based on this analysis we advance the next proposition:

*Proposition 6: There is a set of first best compensation strategies combining fixed and variable pay, contingent and discretionary bonuses, and short-term and long term incentives, such that the relationship between pay and agent motivation is optimized.*

Figure 2 also illustrates a number of other phenomena: below point α motivation falls away rapidly as a result of inequity aversion – effort levels are only restored at point α when the Adams' ratio recovers to an acceptable level; above point β crowding-out means that intrinsic motivation has more or less been eliminated entirely and total motivation has peaked.  If the strong crowding-out conjecture is correct, then at point γ the intrinsic motivation  curve moves from positive to negative and total motivation begins to decline.

It is important to understand what this figure does and does not tell us about executive compensation.  In Figure 2, λ represents the point where total compensation, comprising fixed pay, incentives and rewards, is at its most efficient and effective, and an agent's effort- to-pay ratio is at its highest.  The actual pay of senior executives, which is in practice influenced by other factors such as (often imperfect) labor market conditions, strategic (inter-firm) rivalry, and political (intra-firm) gaming, may in practice be higher.  Executives might be prepared to offer more effort for more incentive pay, but the marginal cost to the employing company of increasing incentive payments may be very high.  This is consistent with the phenomenon of high executive compensation (which may be effective but is not necessarily efficient) and also with proposition 5b, that weak incentives are a more effective and efficient way of motivating agents than strong incentives.

**Corporate Performance**

A complete theory of agency must explain the mechanism which links the job performance of an executive with the performance of the firm. We take as starting point upper echelons theory (Carpenter *et al.*, 2004; Finkelstein *et al.*, 2009; Hambrick & Mason, 1984) which postulates a causal connection between business performance (the dependent variable), the cognitive skills of top managers, their observable personal characteristics (e.g., age, education, experience, socioeconomic background etc.), their strategic choices, and the objective situation (independent variables). We first simplify this a little by taking corporate financial performance to be a function of an agent's performance (as described in the motivation cycle), the performance of other agents, and the external business environment. We then build on the upper echelons approach by postulating a link between the performance of an individual agent $\bar{a}$ (itself a function of his or her ability, motivation and opportunity set), the performance of other agents $\bar{o}$ who, together with agent $\bar{a}$, comprise the top management team, the business strategy (as devised and implemented by the top management team) and the business environment, on the one hand, and business performance on the other hand.

The external business environment is largely outside the control of senior management and hence exogenous to behavioral agency theory. The job performance of other agents, $P_{\bar{o}}$, is endogenous. Indeed, the motivation and performance cycles described in this paper are replicated for all agents fitting the definition of top managers. This generates a final research proposition, that incentive compatibility between agents is a necessary condition of optimal corporate performance. We articulate this as follows:

> *Proposition 7: The incentives and rewards of individual agents must be compatible*
> *with the incentives and rewards of other agents working as part of the same top*
> *management team if firm performance is to be optimized.*

It means, for example, that agents' goals and performance conditions attaching to incentives must be compatible one with another. It also requires inequity aversion to be taken into account within the top management team - there is a strong presumption that individual agents will regard other agents in the same top management team as among their pool of referents for the purposes of equity theory. The desirability of compatible incentives is consistent with the argument that interventions may be necessary in order to align the interest of different members of top management groups (Carpenter *et al.*, 2004; Hambrick, 1994). The incentive compatibility proposition also provides a further argument in favor of weak rather than strong incentives. Roberts (2010) notes that strong incentives may not be appropriate when cooperation between different agents in necessary. Teece, Pisano and Shuen (1997) have pointed out that it is difficult to calibrate individual contributions to a joint effort and have commented that high-powered incentives might well be destructive of cooperative activity and learning.

## DISCUSSION

In her assessment and review of agency theory, Eisenhardt (1989) sets out the main elements of positive agency theory in a table. We repeat this below in Table 2, adding a third column which identifies the areas where behavioral agency theory departs from the standard principal-agent model. According to Eisenhardt the key idea of agency theory is that principal agent relationships should reflect the efficient organization of the costs of information and risk-bearing. The unit of analysis is the contract between principal and agent. The main assumptions are that executives are rational (but see footnote 3), self-interested and risk averse, that there is partial goal conflict between stakeholders, that information is incomplete and not equally shared, and that the overriding organizational objective is efficiency. The problems addressed by the theory involve moral hazard, adverse selection and

how best to share risk, especially where principals and agents have partially differing goals and risk preferences. Proposed solutions to the problems include monitoring through effective corporate governance and outcome-based incentive contracts.

Behavioral agency theory departs from the positive agency framework in three main respects.  First, while positive agency theory focuses on the implications for the firm of costs which arise out of the principal-agent relationship, using efficiency as the main assessment criterion, behavioral agency theory focuses on the relationship between agency costs and performance, using efficiency and effectiveness as the yardsticks.  The objective of an agency contract is to optimize job performance given the constraints of agency costs.  This is achieved at the inflection point $\lambda$ on the agent's proforma pay-effort curve.  Secondly, while agency theory assumes that agents are rational, risk averse and rent seeking, and that there is no non-pecuniary agent motivation, behavioral agency theory proposes a more sophisticated model of man whereby agents are boundedly-rational, loss, risk and uncertainty averse, and where there is a trade-off between intrinsic and extrinsic rewards.  Thirdly, while agency theory assumes a linear relationship between pay and motivation, behavioral agency theory proposes a more complex pay-effort function which is affected by loss, risk and uncertainty aversion, the hyperbolic discounting of deferred rewards, inequity aversion, and the trade-off between intrinsic and extrinsic motivation.

INSERT TABLE 2 ABOUT HERE

**Implications of Behavioral Agency Theory for the Design of Incentives**

Much of the current design thinking about executive compensation ignores behavioral issues and does not take account of agents' preferences, instead falling into the trap of

26

institutional isomorphism (Di Maggio & Powell, 1983), either in the name of "best practice" by following what other firms do (mimetic isomorphism), or by uncritically doing as regulators say (coercive isomorphism).  Behavioral agency theory goes against the current fashion, pointing instead to simpler, more balanced reward systems and more straightforward performance measures.  In particularly, contrary to the logic of agency theory, we argue that high powered incentives are not an efficient and effective way of motivating agents.  It is not possible to construct an incentive contract for an agent or set performance measures which incorporate all the principal's current objectives and are flexible enough to deal with all possible exogenous shocks which might occur during the performance cycle.   Knowledge constraints resulting from an agent's bounded rationality mean that designing very complex incentive contracts in order to tie the principal's and agent's interests as tightly as possible is likely to have an adverse effect of the agent's job satisfaction and work motivation. Furthermore, at high levels of compensation, crowding-out means that intrinsic motivation which is forgone because of an increase in incentives can only be compensated for by proportionately greater increases in extrinsic rewards.  Finally, deferred pay, frequently advocated as a solution to the problem that high levels of executive compensation are seen to be undesirable as a matter of public policy, is in practice an expensive way of paying agents when seen in the context of agent motivation.  These arguments are consistent with the "strength of weak incentives" thesis, described above, as advocated by Roberts (2010).  They contradict the normative arguments of financial economists who advocate  the use of high-powered incentives as a partial remedy for the agency problem (Jensen & Murphy, 1990); see also Bebchuk and Fried (2004: 72).

We argue that, for any group of agents comprising a top management team, there is a balanced set of first best reward strategies, being combinations of fixed and variable pay, contingent and discretionary bonuses, and short and long-term incentives, which allow the

27

relationship between reward costs, agent motivation and job performance to be optimized. In order to maximize firm performance the selected strategy must be compatible with the strategies selected for other agents in the principal's top management team. Identifying these reward strategies is not a simple matter, ideally requiring an understanding of individual differences between agents in terms of their tolerance of risk and inequity and in the way that they discount future rewards. Partly as a result of the complexity involved in designing appropriately simple incentive and reward systems, ex post discretionary payments to agents may sometimes be warranted as partial gift exchanges in the expectation that they will result in reciprocal gifts of effort (Akerlof, 1982).

**Contribution**

Agency theory is a central component of the modern theory of the firm (Jensen, 2000; Roberts, 2004). We have explained that the standard theory of agency has significant shortcomings, especially in its failure in practice to explain the relationship between executive compensation, agent behavior and firm performance. While there is, after Cyert and March (1963 |1992), an extensive literature on the behavioral theory of the firm (see Gavetti, Greve, Levinthal & Ocasio, 2012), we do not yet have a satisfactory behavioral agency theory. This paper takes a significant step in correcting this omission. In particular, it advances a theory of behavioral agency which better explains the mechanisms which connect incentives, agent behavior, and the type of high performance outcomes which shareholders desire. This is an important framework, especially for scholars studying executive compensation.

Positive agency theory, like many theories which have their origins in neoclassical economics, aims to provide accurate predictions about economic phenomena without claiming that its foundational assumptions realistically describe the underlying behavioral

processes (see Friedman, 1953). Wakker (2012), who in turn cites Harré (1970), calls this paramorphism. This approach is self-evidently flawed when neither the predictions nor the underlying processes match reality. Behavioral agency theory, on the other hand, aims to explain economic phenomena by reference to descriptions of underlying processes which do match reality. Wakker (2012) calls this homeomorphism. Behavioral economists argue that homeomorphism is more likely to generate useful explanations of actual economic phenomena and hence is a better approach to theory building.

Part of the validity which we claim for behavioral agency theory is based on the way in which it adapts and integrates existing theory. Steel and König have emphasized the importance of consilience in theory development, arguing that: "if a theory can be shown to have consilience, its scientific validity is vastly improved, since it represents different avenues of inquiry coming to similar conclusions" (Steel & Konig, 2006: 889). A major contribution of this paper lies in the way in which it integrates a number of different literatures: in particular, the neoclassical economic theory of agency (after Jensen & Meckling, 1976); work motivation theory (for example Locke & Latham, 1984; Steel & Konig, 2006; and Vroom, 1964); the literature on choices, values, heuristics and biases (after Kahneman & Tversky, 1979; and Tversky & Kahneman, 1992); and the upper echelons approach to strategic leadership (after Hambrick & Mason, 1984). Our paper also makes a contribution to the management literature by updating Eisenhardt's (1989) review of agency theory for management scholars, incorporating new ideas from behavioral economics. In addition, and significantly, the paper has important implications for practice in the way that it advocates the use of balanced executive reward strategies and weak incentives.

29

**Conclusion**

Formally, behavioral agency theory comprises four inter-connected equations, two figures, ten propositions, and a supporting narrative. Equation (1), after Boxall and Purcell (2003), connects an agent's job performance with his or her ability, motivation and opportunity set. Equation (2), after Vroom (1964), which is in turn modified by the inclusion of time discounting, risk and loss aversion to become Equation (3), after Steel and König (2006), explains the relationship between compensation and agent motivation. Equation (4) distinguishes between intrinsic and extrinsic motivation and identifies a potential trade-off between the two. Figure 1 explains the place of agent performance and work motivation in the firm's performance cycle and Figure 2 illustrates an agent's pay-effort curve.

In this paper we have sought to provide a better understanding of the micro-foundations of agency theory, especially as it applies to executive compensation, based on a more realistic set of assumptions about agent behavior. We hope that others will join us in further developing behavioral agency theory, in testing it empirically, and in identifying other implications for business practice.

# REFERENCES

Adams, J. S. 1965. Inequity in social exchange. In L. Berkowitz (Ed.), *Advances in experimental social psychology*: Academic Press, New York.

Ainslie, G. 1991. Derivation of 'Rational' Economic Behavior from Hyperbolic Curves. *American Economic Review, 81(2)*: 334-340.

Ainslie, G., & Haslam, N. 1992. Hyperbolic discounting. In G. Loewenstein and N. Haslam (Eds.), *Choices over time*. New York: Russell Sage Foundation.

Akerlof, G. 1982. Labor contracts as partial gift exchange. *Quarterly Journal of Economics, 97*: 543-569.

Alchian, A., & Demsetz, H. 1972. Production, information costs and economic organization. *American Economic Review, 62*: 777-795.

Arrow, K. 1985. The economics of agency. In J. Pratt and R. Zeckhauser (Eds.), *Principals and agents: the structure of business*. Cambridge, MA: Harvard Business School Press.

Balkin, D. 2008. Explaining High US CEO Pay in a Global Context: An Institutional Approach. In L. Gomez-Mejia and S. Werner (Eds.), *Global Compensation: Foundations and Perspectives*. London: Routledge.

Barnard, C. 1968. *The functions of the executive*. Cambridge, MA: Harvard University Press.

Bebchuk, L., & Fried, J. 2004. *Pay without performance – the unfilled promise of executive compensation*. Cambridge, Massachusetts: Harvard University Press.

Bebchuk, L., Fried, J., & Walker, D. 2002. Managerial power and rent extraction in the design of executive compensation. *University of Chicago Law Review, 69*: 751 – 846.

31

Becker, G. 1993. *Human capital. A theoretical and empirical analysis with special reference to education* (Third ed.). Chicago: University of Chicago Press.

Besley, T., & Ghatak, M. 2005. Competition and incentives with motivated agents. *American Economic Review, 95*: 616-636.

Boivie, S., Lange, D., McDonald, M., & Westphal, J. 2011. Me or we: the effects of CEO organizational identification on agency costs. *Academy of Management Journal, 54(3)*: 551-576.

Boxall, P., & Purcell, J. 2003. *Strategy and human resource management*. Basingstoke, Hants: Palgrave Macmillan.

Bratton, W. 2005. The Academic tournament over executive compensation. *California Law Review, 93(5)*.

Camerer, C., Loewenstein, G., & Rabin, M. 2004. *Advances in Behavioral Economics*. Princeton, NJ: Princeton University Press.

Carpenter, M., Geletkanycz, M., & Sanders, W. 2004. Upper echelons research revisited: antecedents, elements, and consequences of top management team composition. *Journal of Management, 30(6)*.

Charreaux, G. 2002. Positive agency theory: place and contributions. In E. Brousseau and J. Glachant (Eds.), *The economics of contracts*. Cambridge, England: Cambridge University Press.

Combs, J., & Skills, M. 2003. Managerialist and Human Capital Explanations for Key Executive Pay Premiums. *Academy of Management Journal, 46(1)*: 63-73.

32

Cyert, R., & March, J. 1963 |1992. A behavioral theory of the firm. Cambridge, MA: Blackwell.

Deci, E., & Ryan, R. 1985. *Intrinsic motivation and self-determination in human behavior*. New York: Plenum.

Devers, C., Cannella, A., Reilly, G., & Yoder, M. 2007. Executive compensation: a multidisciplinary review of recent developments. *Journal of Management, 33*: 1016-1072.

Di Maggio, P., & Powell, W. 1983. The iron cage revisited: institutional isomorphism and collective rationality in organizational fields. *American Sociological Review, 48*: 147-160.

Eisenhardt, K. M. 1989. Agency theory: an assessment and review. *Academy of Management Review, 14*: 57-74.

Fehr, E., & Schmidt, K. 1999. A theory of fairness, competition, and cooperation. *The Quarterly Journal of Economics, 114*: 817-868.

Finkelstein, S., Hambrick, D., & Cannella, A. 2009. *Strategic leadership: theory and research on executives, top management teams, and boards*. Oxford: Oxford University Press.

Foss, N. 2010. Bounded rationality and organizational economics. In P. Klein and M. Sykuta (Eds.), *The Elgar companion to transaction cost economics*. Cheltenham: Edward Elgar Publishing Limited.

Frederick, S., Loewenstein, G., & O'Donoghue, T. 2002. Time Discounting and Time Preference: A Critical Review. *Journal of Economic Literature, 40(2)*: 351-401.

Frey, B. S., & Jegen, R. 2001. Motivation crowding theory. *Journal of Economic Surveys, 15*: 589-611.

Friedman, M. 1953. The Methodology of Positive Economics. In M. Freidman (Ed.), *Essays in Positive Economics*. Chicago: University of Chicago Press.

Frydman, C., & Jenter, D. 2010. CEO compensation. *Annual Review of Financial Economics, 2*: 75-102.

Gabaix, X., & Landier, A. 2008. Why has executive pay increased so much? *Quarterly Journal of Economics, 123*: 49-100.

Gavetti, G., Greve, H., Levinthal, D., & Ocasio, W. 2012. The behavioral theory of the firm: assessment and prospects. *The academy of management annals*. New York: Routledge Taylor & Francis Group.

Gomez-Mejia, L., Berrone, P., & Franco-Santos, M. 2010. *Compensation and organizational performance - theory, research and practice*. Armonk, NY: M.E. Sharpe, Inc.

Gomez-Mejia, L., & Wiseman, R. 1997. Reframing executive compensation: an assessment and outlook. *Journal of Management, 23*: 291-374.

Graves, S., & Ringuest, J. 2012. Patient Decision Making: Exponential versus Hyperbolic Discounting. *Managerial and Decision Economics*: Wiley Online Library.

Grossman, S., & Hart, O. 1983. An Analysis of the Principal-Agent Problem. *Econometrica, 51(1)*: 7-45.

Hambrick, D. 1994. Top management groups: a conceptual integration and reconsideration of the "team" label. In B. Staw and L. Cummings (Eds.), *Research in Organizational Behavior*: 171-213. Beverly Hill: JAI Press.

Hambrick, D., & Mason, P. 1984. Upper echelons: the organization as a reflection of its top managers. *Academy of Management Review, 9(2)*: 193-206.

Harre, R. 1970. *The Principles of Scientific Thinking*. London: Macmillan.

Hart, O. 1995. *Firms, Contracts and Financial Structure*. Oxford: Oxford University Press.

Jensen, M. 1983. Organization theory and methodology. *The Accounting Review, 58*: 319-339.

Jensen, M. 1998. *Foundations of organizational strategy*. Cambridge, MA: Harvard University Press.

Jensen, M. 2000. *A theory of the firm: governance, residual claims, and organizational forms*. Cambridge, MA: Harvard University Press.

Jensen, M., & Meckling, W. 1976. Theory of the firm: managerial behaviour, agency costs and ownership structure. *Journal of Financial Economics, 3*: 305-360.

Jensen, M., & Meckling, W. 1994. The nature of man. *Journal of Applied Corporate Finance, 7(2)*.

Jensen, M., & Murphy, K. 1990. Performance pay and top-management incentives. *Journal of Political Economy, 98*: 225-264.

Kahneman, D., & Tversky, A. 1979. Prospect theory – an analysis of decision under risk. *Econometrica, 47*: 263-291.

Kreps, D. 1997. Intrinsic motivation and extrinsic incentives. *The American Economic Review, 87*: 359-364.

Lazear, E. P., & Rosen, S. 1981. Rank-Order Tournaments as Optimum Labor Contracts. *Journal of Political Economy, 89*: 841 – 864.

Leibenstein, H. 1966. Allocative efficiency vs. x-efficiency. *The American Economic Review, 56*: 392-415.

Lepper, M., & Greene, D. 1978. Introduction. In M. Lepper and D. Greene (Eds.), *The hidden costs of reward - new perspectives on the psychology of human motivation.* Hillsdale, NJ: Lawrence Erlbaum Associates.

Lewin, K. 1938. The conceptual representation and measurement of psychological forces. *Contributions to Psychological Theory, 4*: 247.

Locke, E., & Latham, G. 1984. *Goal setting - a motivational technique that works*. New Jersey: Prentice-Hall Inc.

Locke, E., & Latham, G. 1990. *A theory of goal setting and task performance*. Englewood Cliffs, NJ: Prentice Hall.

Locke, E., & Latham, G. 2002. Building a practically useful theory of goal setting and task motivation - a 35 year odyssey. *American Psychologist, 57*: 705-717.

Michelman, F. R. 1967. Property, utility, and fairness: comments on the ethical foundation of 'just compensation' law. *Harvard Law Review, 80*: 1165-1257.

Mintzberg, H. 1997. *The nature of managerial work*. London: Prentice Hall.

Mintzberg, H. 2009. *Managing*. London: FT Prentice Hall.

36

Pepper, A. 2006. *Senior executive reward – key models and practices*. Aldershot: Gower.

Pratt, J., & Zeckhauser, R. 1985. Principals and agents: an overview. In J. Pratt and R. Zeckhauser (Eds.), *Principals and agents: the structure of business*. Cambridge, MA: Harvard Business School Press.

Rebitzer, J., & Taylor, L. 2011. Extrinsic rewards and intrinsic motives: standard and behavioral approaches to agency and labor markets. In O. Ashenfelter and D. Card (Eds.), *Handbook of labor economics*: 701-772. Amsterdam: North-Holland.

Roberts, J. 2004. *The modern firm*. Oxford: Oxford University Press.

Roberts, J. 2010. Designing incentives in organizations. *Journal of Institutional Economics, 6*: 125-132.

Roberts, J. 2011. Weak incentives: when and why. In D. Marsden (Ed.), *Employment in the lean years*: 185-198. Oxford: Oxford University Press.

Ross Ashby, W. 1956 |1976. *Introduction to cybernetics*. London: Methuen.

Ross, S. 1973. The economic theory of agency: the principal's problem. *American Economic Review, 63*: 134-139.

Sanders, G., & Carpenter, M. 2003. A behavioral agency theory perspective on stock repurchase program announcements. *Academy of Management Journal, 46(3)*: 160-178.

Savage, L. 1954. *The foundations of statistics*. New York: Wiley.

Shafir, E., Diamond, P., & Tversky, A. 1997. Money illusion. *The Quarterly Journal of Economics, 112*: 341-371.

Simon, H. 1945 |1997. *Administrative behavior* (Fourth ed.). New York: The Free Press.

Simon, H. 1957 |1982. The compensation of executives; first published in Sociometry (1957) Vol. 20. *Models of bounded rationality, volume 2: behavioral economics and business organisation* 47-48. Cambridge, MA: The MIT Press.

Sliwka, D. 2007. Trust as a signal of a social norm and the hidden costs of incentive schemes. *American Economic Review, 97(3)*: 999-1012.

Spence, M. 1973. Job market signalling. *Quarterly Journal of Economics, 87(355-374)*.

Spence, M., & Zeckhauser, R. 1971. Insurance, information and individual action. *American Economic Review, 61*: 380-387.

Steel, P., & Konig, C. 2006. Integrating theories of motivation. *Academy of Management Review, 31*: 889 – 913.

Teece, D., Pisano, G., & Shuen, A. 1997. Dynamic capabilities and strategic management. *Strategic Management Journal, 18(7)*.

Tosi, H., & Gomez-Mejia, L. 1989. The decoupling of CEO pay and performance: an agency theory perspective. *Administrative Science Quarterly, 34*: 169-189.

Tosi, H., Werner, S., Katz, J., & Gomez-Mejia, L. 2000. How much does performance matter? A meta-analysis of CEO pay studies. *Journal of Management, 26*: 301-339.

Tversky, A., & Kahneman, D. 1992. Advances in Prospect Theory: Cumulative Representation of Uncertainty. *Journal of Risk and Uncertainty, 5*: 297-323.

Ungson, G., & Steers, R. 1984. Motivation and Politics in Executive Compensation. *Academy of Management Review, 9(2)*: 313-323.

Vroom, V. H. 1964. *Work and motivation*. New York: Wiley.

Wade, J., O'Reilly, C., & Pollock, T. 2006. Overpaid CEOs and Underpaid Managers: Fairness and Executive Compensation. *Organizational Science, 17(5)*: 527-544.

Wakker, P. 2012. *Prospect theory*. Cambridge: Cambridge University Press.

Wiseman, R., & Gomez-Mejia, L. 1998. A behavioral agency model of managerial risk taking. *Academy of Management Review, 23*: 133-153.

# FOOTNOTES

[1] Some commentators (e.g., Roberts, 2011) imply that Jensen and Murphy's empirical evidence is not contrary to agency theory, but suggests instead that it means the (normative) recommendations of agency theory have not been followed in practice. An argument in this form, implying that the absence of two factors (incentive pay and high performance) can be interpreted as evidence of a causal connection between the two phenomena (so that more of the first factor will necessarily lead to more of the second) is hardly justified. It also appears to confuse the positive theory of agency (which should be capable of explaining the world as it is) with normative theory. In practice, as argued long ago by Herbert Simon in 1957 and demonstrated empirically by Gabaix and Landier in 2008, CEO pay is much more closely correlated with company size than company performance.

[2] Jensen (1983) identifies two different strands in the literature on agency theory. He calls these the "positive theory of agency" and the "principal-agent" literature. Eisenhardt (1989) describes the latter as a "general theory of the principal-agent relationship", while Wiseman and Gomez-Mejia (1998) call it "normative agency theory". Positive agency theory focuses on the special case of the principal-agent relationship between owners and managers of large corporations (Charreaux, 2002; Eisenhardt, 1989; Jensen, 1983). Normative agency theory aims to provide a formal theory of the principal-agent relationship in all its guises, including employer-employee, lawyer-client, buyer-supplier etc., (Eisenhardt, 1989).

[3] Eisenhardt (1989) states that positive agency theory also assumes bounded rationality, but we can find no other reference to this in the agency theory literature. After the first, formative papers on agency theory, Jensen and Meckling (1994) later develop the resourceful, evaluative, maximizing model of man (REMM) which they say is consistent with agency theory, but this is still a rational choice model. They subsequently develop a second framework, the pain avoidance model (PAM), but they do not seek to integrate this into agency theory.

[4] Christen et al., (2006) point out that motivation (wanting to work hard) is not the same as effort (working hard and, in doing so, expending time and energy). However, in much the same way that revealed preference is a marker of mental preference, so effort can be thought of as a marker of motivated behavior (see Martin and Tesser, 2009). This means that, in the absence of coercion, effort can be taken to imply the presence of motivation even if motivation does not necessarily result in the expenditure of effort.

[5] Jensen and Meckling do not explicitly mention expenditure on incentives and rewards, i.e., the actual cash costs of incentivizing and rewarding agents, although such expenditure would seem self-evidently to be part of

the cost of agency. Incentive and reward costs can be further broken down into the costs of providing incentives and rewards in the optimal form and mix plus any additional costs incurred in incentivizing and rewarding agents in a way which is sub-optimal. In order to be precise we use the following terminology in this article: total compensation or pay ($\omega$) is the sum of fixed pay and variable pay; variable pay is itself the sum of incentives (awarded ex ante) and discretionary rewards (awarded ex post). "Compensation" and "pay" are treated as synonyms. Use of the sign "$\omega$" follows the convention in labor economics of taking "$\omega$" as the symbol for wages.

[6] There are many different definitions of "bounded rationality". We follow Williamson, who explains that rationality is subject to neuro-physiological rate and storage limits on the powers of agents to receive, store, retrieve, and process information without error (after Williamson, 1975, p.21). Williamson also talks about a further element of bounded rationality, which he calls "language limits", being the constraints on individuals to communicate comprehensively in such a way that they are fully understood by others, but this element is not really relevant to the current article. Foss (2010) provides an elegant summary description of bounded rationality, which he describes in terms of (1) limitations in the human capacity to process information; (2) attempts to economize on mental effort by relying on short-cuts or heuristics; and (3) a consequence of the fact that cognition and judgement are subject to a wide range of biases and errors.

[7] In this paper we largely ignore the Knightian distinction between risk (probabilistic outcomes) and uncertainty (indeterminate outcomes), instead treating "risk" and "uncertainty" as synonymous.

[8] Note however that Barnard used the term "efficiency" in an entirely different sense: to Barnard an organization is "efficient" if it satisfies the motives of its members.

[9] Expectancy, a measure of probability, takes values between 0 and +1. Instrumentality takes values between +1 (meaning it is believed that the first outcome will certainly lead to the second outcome) and -1 (meaning it is believed that the second outcome is impossible in the event of the first outcome).

[10] Temporal motivation theory combines expectancy and instrumentality into one operator, which Steel and König call "expectancy" but which is essentially the same thing as subjective probability after Savage (1954). While this loses some of the richness of Vroom's conceptualization of expectancy and instrumentality (especially the possibility that instrumentality may be negative) it is a pragmatic simplification of the theory and hence is followed here.

41

[11] A more complex way of representing Steel and König's motivation function, which distinguishes between gains and losses and hence accounts for loss aversion is:

$$X_i = \sum \left\{ \frac{E_{ik}^{+} \times V_{k}^{+}}{1 + \delta^{+}t} \right\} - \sum \left\{ \frac{E_{ik}^{-} \times V_{k}^{-}}{1 + \delta^{-}t} \right\}$$

This expression of the formula explicitly recognizes that the expectancy, valence and the average discount factor will differ for gains (represented by $^{+}$) and losses (represented by $^{-}$).

**TABLE 1: Assumptions about the Nature of Man under Positive Agency Theory and**

**Behavioral Agency Theory**

| Assumption | Economic man | Behavioral economic man |
|---|---|---|
| Principal's risk preference | Principals are risk neutral | As for agency theory |
| Agent's utility function | Agents are rent seeking, agent's utility is positively contingent on pecuniary incentives and negatively contingent on effort | As for agency theory, but subject to constraints relating to rationality, motivation, loss, risk, uncertainty and time preferences |
| Agent's rationality | Agents are rational | Agents are boundedly rational, i.e., subject to neuro-physiological rate and storage limits on the powers of agents to receive, store, retrieve, and process information without error |
| Agent's motivation | There is no non-pecuniary agent motivation | Motivation is both intrinsic and extrinsic. Intrinsic and extrinsic motivation are neither independent nor additive. |
| Agent's risk preference | Agents are risk averse | Agents are loss averse below a gain/loss inflection point; otherwise risk averse |
| Agent's time preferences | Agents' time preferences are calculated according to an exponential discount factor | Agents' time preferences are calculated according to a hyperbolic discount factor |
| Agent's preference for perceived equitable pay | Not defined | Agents are inequity averse |

43

**TABLE 2: Overview of Positive Agency Theory vs. Behavioral Agency Theory**

|  | **Positive agency theory** (after Eisenhardt, 1989) | **Behavioral agency theory** |
|---|---|---|
| Key idea | The primary importance of aligning the interests of principals and agents. The principal-agent relationship should reflect efficient management of the costs of information and risk-bearing | The primary importance of agent performance and work motivation. The principal-agent relationship should reflect the efficient *and effective* management of the relationship between executive compensation, firm performance and shareholder interests |
| Unit of analysis | Contract between principal and agent | Contract between principal and agent |
| Human assumptions | Agents are rational, self-interested, risk averse | Agents are boundedly rational, loss, risk and uncertainty averse, hyperbolic time discounters, inequity averse, and there is a trade-off between intrinsic and extrinsic motivation |
| Organizational assumption | Partial goal conflict between principals and agents, efficiency as the main performance criterion, information asymmetry | Partial goal conflict between principals and agents, efficiency *and effectiveness* as the main performance criterion, information asymmetry |
| Information assumption | Asymmetric information and incomplete contracting | As for agency theory; goal setting used as a pragmatic solution to information asymmetry |
| Primary factor(s) determining the principal-agent relationship | The principal's wish to align the agent's objectives with the principal's owns objectives (alignment) | The principal's wish to align the agent's objectives with the principal's own objectives (alignment) and to motivate agents to give high performance, given their abilities and opportunities (motivation) |
| Contracting problems | Moral hazard and adverse selection | As for agency theory |

44

| | | |
|---|---|---|
| Key mechanisms | Monitoring and incentive contracts | As for agency theory, except that incentive contracts can also help to meet the motivation objective |
| Problem domain | Where principals and agents have different goals and risk preferences e.g., regulation, compensation, vertical integration, transfer pricing | As for agency theory; especially relevant to executives and executive compensation |

45

**FIGURE 1**
**Agent's Job Performance and Work Motivation Cycle**

Intrinsic motivation

Loss, risk & uncertainty
aversion

Time
discounting

P2a, b

P3

P1a, b

Extrinsic
motivation
(i.e., incentives)

Agent's motivation

Agent's
job
performance

Rewards
P6, P7

P4

P5a, b

Inequity aversion

Goal setting, contracting
& monitoring

Feedback

46

**FIGURE 2**
**Agent's Pay-Effort curve**



47

489

# Executive Compensation: Where We Are, and How We Got There

**Kevin J. Murphy**
University of Southern California, Marshall School of Business, Los Angeles, CA, USA

## Contents

# 1.  INTRODUCTION

The first decade of the new century brought significant changes to executive compensation in large US companies. Rocked by scandals ranging from accounting fraud to option backdating—coupled with suspicions that Wall Street bonuses led to excessive risk taking that triggered the financial crisis—compensation committees faced a plethora of new pay-related laws and tax, accounting, and disclosure rules designed to stem perceived abuses in executive pay. After more than tripling (after inflation) during the 1990s stock-option explosion, the median total pay for chief executive officers (CEOs) in the S&P 500 remained relatively stagnant in the early 2000s, and indeed

even declined during the 2008–2009 Great Recession. But the flattening of pay levels belied significant structural changes in the composition of pay, as companies adapted to the new regulations and jettisoned stock options in favor of restricted stock. Moreover, realized pay for top-level executives was postured for a new explosion in the second decade of the 2000s, as stock and options granted near the bottom of the market in 2009 became vested and exercisable. These trends suggest the outrage over executive pay—recently reflected by the "Occupy Wall Street" movements and in calls from the Obama administration for increased tax rates for "millionaires and billionaires"—will likely continue unabated over the next several years.

The recent controversies over executive pay are not the first—nor will they be the last—time that executive compensation has sparked outrage and calls for regulation and reform. Indeed, scrutinizing, criticizing, and regulating high levels of executive pay has been an American pastime for nearly a century. In 1932, for example, controversies surrounding high salaries for executives in bailed-out railroads led to pay disclosures and pay caps; disclosure requirements were soon extended to banks, utilities, and large corporations, and further extended to all publicly traded companies following the 1933 and 1934 Securities Acts. Outrage over perceived excesses in "restricted stock option plans" in the 1960s led Congress to prohibit repricing, reduce maximum expiration terms, restrict exercise prices, and extend required holding periods after exercises. In the 1980s, Congress imposed large tax penalties on firms paying (and executives receiving) large severance payments following a change in control, and in the 1990s non-performance-based pay exceeding $1 million was deemed unreasonable and therefore not deductible as an ordinary business expense for corporate income tax purposes. Therefore, the recent backlash over executive pay associated with the accounting and backdating scandals and the financial crisis—triggering Sarbanes-Oxley, new disclosure and accounting rules, restrictions on deferred compensation, and myriad pay regulations under the Dodd–Frank Act—continues a tradition of regulatory responses to perceived excesses and abuses in top-level pay.

The purpose of this study is to document the current state of executive compensation and to show how the level and structure of CEO pay over the past century has evolved in response to economic, institutional, and political factors. My intention is not to provide a comprehensive survey of the academic literature on executive compensation (or even a systematic update of Murphy (1999)), but rather to document a body of facts to guide future theoretical and empirical research in the area. I show that government intervention into executive compensation—largely ignored by researchers—has been both a response to and a major driver of time trends in CEO pay. There have been two broad patterns for government intervention into CEO pay. The first pattern is aptly described as knee-jerk reactions to isolated perceived abuses in pay, leading to disproportionate "one-size-fits-all" responses and a host of unintended and undesirable consequences. The second pattern—best described as "populist" or "class

warfare"—arises in situations where CEOs (and other top executives) are perceived to be getting richer when lower-level workers are suffering. Beyond these two broad patterns, indirect intervention in the form of accounting rules, securities laws, broad tax policies, and listing requirements have also had direct impact on the level and composition of CEO pay. In most cases, companies and their executives have responded to the interventions by circumventing or adapting to the reforms, usually in ways that increased pay levels and produced other unintended (and typically unproductive) consequences.

More broadly, this study provides institutional context useful in "explaining" time trends in the level and structure of CEO pay. As emphasized by Frydman and Jenter (2010) and explored below in Section 5, the academic literature focused on explaining these trends is roughly divided into two camps: the "efficient contracting" camp and the "managerial power" camp. The efficient-contracting camp—rooted in optimal contracting theory—maintains that the observed level and composition of compensation reflects a competitive equilibrium in the market for managerial talent, and that incentives are structured to optimize firm value. The managerial-power camp—exemplified in a series of papers by David Yermack, Lucian Bebchuk and Jesse Fried—maintains that both the level and composition of pay are determined not by competitive market forces but rather by captive board members catering to rent-seeking, entrenched CEOs. Frydman and Jenter (2010) conclude that neither camp offers convincing explanations for cross-sectional and time-series patterns in the data.

The efficient-contracting and managerial power camps are not mutually exclusive. For example, in a series of papers designed to explain the escalation in option grants in the 1990s, I have argued that stock options were granted in such large quantities to so many employees in the 1990s because boards and executives (erroneously) perceived options to be essentially free to grant.[1] This explanation might be viewed as a combination of both camps: directors yielded to shareholder pressure to tie more closely to equity values, but were duped by managers into the idea that options were free to grant, thus leading to massive grants without any noticeable reductions in other forms of pay. However, as will become clear in Section 3.7 below, a more complete explanation must include the role of government: the option explosion in large part caused by changes to tax and accounting rules coupled with changes in disclosure, holding, and listing requirements.

In essence, the efficient-contracting camp views executive pay as mitigating agency problems between executives and shareholders, while the managerial-power camp views excessive pay as symptomatic of agency problems between shareholders and board members (who often own only a trivial fraction of their firm's common stock and who are in no sense perfect agents for the shareholders who elected them). The reason government

---

[1] See, for example, Murphy (2002, 2003) and Hall and Murphy (2003).

intervention into executive pay adds an important new dimension to the analysis is because the interests of the government differ significantly from those of shareholders, directors, or executives. In particular, as will become evident from the legislative history in Section 3 below, Congressional (and, more generally, public) outrage over executive pay is almost always triggered by perceived excesses in the *level* of compensation without regard to incentives and company performance, and the regulatory responses have also fixated on pay levels (albeit with little effect).

Limitations on the form of government intervention add another interesting dimension to the agency problem. In most circumstances, Congress has stopped short of directly capping the level of pay or imposing restrictions on its structure.[2] However, Congress controls the tax code (including individual and corporate tax rates, punitive excise taxes, and defining what compensation is "reasonable" and therefore deductible by the company), and has routinely used tax rules to regulate pay. In addition, Congress (through its influence on the SEC) indirectly controls disclosure requirements, long the favorite (and singularly most ineffective) tool used to control perceived abuses in pay. Ultimately, attempts to regulate the level of pay through tax and disclosure rules (instead of direct pay caps) have allowed plenty of scope for circumvention and opportunism and other unintended consequences, often leading to the next round of scandals and government responses.

Section 2 ("Where We Are") analyzes the level and structure of CEO pay packages, discusses measurement issues, explores 1970–2011 time trends and, more generally, serves as a primer on executive compensation. I distinguish between three different measures of total compensation: (1) grant-date pay (based on grant-date values for stock and options, and target values for bonuses); (2) realized pay (based on the vesting of stock awards and the gains from exercising options); and (3) risk-adjusted pay (expected pay from the perspective of risk-averse CEOs). I document the dramatic increase in CEO pay during the 1990s, driven primarily by an unparalleled escalation of stock option grants, and the flattening of pay during the early 2000s (as firms replaced option grants with stock awards). In addition, I provide 1992–2011 time-series evidence on the relation between CEO wealth and shareholder wealth and stock-price volatilities, and discuss incentive issues related to bonus plans and earnings announcements.

---

[2] Congress has occasionally attempted to cap wage increases. For example, the World War II Stabilization Act of 1942 froze wages and salaries (for all workers, not just executives), and the 1971 Nixon wage-and-price controls imposed a 5.5% limit on increases in executive pay (the limit being binding for company-defined groups of executives, but not necessarily for individual executives). In addition, Congress has occasionally imposed restrictions on individual pay components, such as Sarbanes-Oxley's prohibition on company-provided loans. More recently, Congress directly (and enthusiastically) regulated both the level and structure of pay for executives in financial services firms receiving assistance under Treasury's Troubled Asset Relief Program ("TARP"), see Section 3.8.5 below.

Section 3 ("How We Got There") provides a history of CEO pay in the United States, emphasizing the causes and consequences of government interventions, which have substantially prohibited what would otherwise be highly desirable and productive pay practices. I begin by examining the controversies leading to the first public disclosures of executive pay in the 1930s, which in turn laid the groundwork for all future controversies of, and interventions into, US CEO pay. I document the rise and fall of restricted stock options in the 1950s, created and ultimately destroyed by changes in tax rules. I discuss how wage-and-price controls and a stagnant stock market facilitated an explosion in perquisites in the 1970s; the surrounding controversy led to new tax and disclosure rules (but did not seem to lead to a reduction in perquisites). I show how penalties on golden parachutes in the 1980s appear to have increased the prevalence of change-in control plans; tax gross-ups, early exercise of stock options, and employment agreements. While the increase in option grants in the 1990s in part reflected increased pressure from shareholders to tie CEO pay more closely to performance, I show that the option explosion is largely attributed to tax, accounting and disclosure rules coupled with changes in holding and listing requirements that favored stock options over other forms of incentive compensation. Next, I speculate that the increased reliance on options helped fuel the accounting and backdating scandals in the early 2000s, which in turn led to a variety of government responses and subsequent changes in compensation (including the move towards restricted stock). I then discuss the pay restrictions for recipients of government bailouts during the financial crisis. Finally, I discuss the ongoing implementation of the Dodd–Frank Act.

Section 4 provides international comparisons of CEO pay, based largely on my joint work with Nuno Fernandes, Miguel A. Ferreira, and Pedro Matos (Fernandes et al., 2012). Based on recently available data from 14 countries with mandatory pay disclosures—we show that the stylized fact that US CEOs earn substantially more than foreign CEOs is wrong, or at least outdated. In particular, the "US Pay Premium" became statistically insignificant by 2007 and largely reflects a risk premium for stock-option compensation (which remains more prevalent in the United States than in other countries). In reaching this conclusion, we control not only for the "usual" firm-specific characteristics (e.g. industry, firm size, volatility, and performance) but also for governance characteristics that systematically differ across countries. The remaining differences in pay are largely explained by evolutionary differences in the politics of pay. In particular, Section 3 showed that CEO pay reflects, in part, political responses to perceived (or actual) abuses in pay. Since those perceived abuses differ across countries, the evolution of pay has also differed. For example, CEO pay became highly controversial in both the United States and the United Kingdom in the early 1990s. In the United States, the (likely unintended) result of the controversy was the explosion in stock option grants. In the United Kingdom, the result of a slightly different controversy was

to essentially move away from options in favor of performance shares and other forms of equity-based compensation.

Section 5 uses the results in the prior sections to suggest a general theory of executive compensation. I argue that viewing efficient contracting and managerial power as competing hypotheses to "explain" executive compensation has not been productive, since the hypotheses are not mutually exclusive and because they ignore critical political factors and other influences on pay. Ultimately, what makes CEO pay interesting, complicated, and worthy of continued investigation is that the paradigms co-exist and interact.

## 2. WHERE WE ARE: A PRIMER ON EXECUTIVE COMPENSATION

### 2.1 Measuring Executive Pay

Underlying every intra-firm, cross-sectional, cross-country, or time-series analysis of executive compensation is an assumption (too often implicit) about how to measure the total compensation received by the executives. If executives were simply paid a base salary set at the beginning of each year, it would be easy to compare salaries across executives (within a firm or across firms, industries, and countries) to identify the highest paid, to compare salaries across years to determine how pay has changed over time, and to compare executive salaries to wages paid in other occupations. But consider the following:

- Executives receive compensation in a dizzying array of forms, including base salaries, annual bonuses, long-term incentives, restricted stock, performance shares (i.e., restricted stock with performance-based vesting), stock options, retirement benefits, and perquisites ranging from health benefits to club memberships and personal use of the corporate jet.
- Many of these forms of compensation depend on performance measured over a single or multiple years, and it is not obvious how (or when) to measure them. For example, stock options (which give the executive the right, but not the obligation, to buy a share of stock at a predetermined price) typically have terms of up to ten years. Should stock options be "counted" as compensation when granted, or only when exercised?
- In addition, executives routinely receive lump-sum amounts at various points in time, such as signing bonuses when joining their firms, severance payments upon termination, and change-of-control payments when their companies are taken over. Moreover, some payments "earned" while employed (such as defined-benefit pension obligations) are not paid until long after the executive is retired and his compensation is no longer reported (or sometimes paid as a lump-sum upon retirement). Again, it is not obvious how, or when, to measure these aspects of compensation.

- Finally, different components of compensation impose different amounts of risk on executives. The payoffs from stock options, for example, are inherently more risky than are payoffs from restricted stock, which in turn are more risky than base salaries. Risk-averse and undiversified executives will naturally place a lower value on riskier forms of compensation, and yet most studies of executive pay simply (and blindly) add together these different forms of compensation. The "risk premia" that executives attach to different forms of compensation depend on unobservable characteristics such as risk aversion and diversification, and it is not obvious how to add or how to weight the various components.

### 2.1.1  "Grant-Date" vs. "Realized" Pay

While the ultimate value of stock awards and stock options is not known until the stock vests and the options are exercised, these equity awards clearly have a value upon grant. Perhaps the most critical choice facing researchers in executive compensation is whether to measure the compensation associated with equity awards as the amount actually realized upon vesting and exercise, or to assign an "ex ante" grant–date value. Most academic research on executive compensation since the mid-1980s has adopted the ex ante approach, valuing stock awards as the fair market value on the date of grant (i.e. the grant-date stock price multiplied by the number of shares granted), and valuing stock options on the grant date using some variant of the Black and Scholes (1973) formula.

When total compensation is measured using grant-date values, it is routinely referred to as *expected compensation* to distinguish it from *realized compensation* as measured at the time the stock vests and the options are exercised.[3] However, calling the grant-date pay "expected" is somewhat loose:

- For restricted shares (i.e. shares to be delivered at a future point in time), the grant-date stock price is the discounted expected value only if there are no performance hurdles, no dividends (or if the executive receives dividends on restricted shares, which is common) and only if there is no risk of forfeiture (i.e. no risk that the employment relation is terminated by either party prior to vesting).
- For stock options, the Black–Scholes value is the discounted expected payoff of a non-forfeitable European option for an executive who can perfectly hedge away the risk of the option (or, alternatively, the expected payoff under the risk-neutral distribution discounted at the risk-free rate).
- As discussed below in Section 2.1.2, the grant–date value (for either stock or option awards) is not a measure of value from the perspective of risk-averse undiversified

---

[3] Standard & Poor's ExecuComp database—the most widely used data in executive compensation research—defines grant-date and realized compensation as "TDC1" and "TDC2", respectively. However, since the value of restricted shares upon vesting has only been disclosed since 2006, ExecuComp actually measures TDC2 using grant-date values for restricted shares (and exercise gains for options).

executives who cannot hedge away the risk. However, with appropriate adjustments for dividends, forfeiture, dilution, and (for options) early exercise, the grant-date value can be an appropriate estimate of the cost to the company of granting restricted stock or options.

Similarly, bonus plans have a "grant-date value" typically measured as the target bonus, paid when the company achieves (usually accounting-based) target performance. However, even when target performance equals expected performance, the target bonus is only the "expected bonus" when the rewards and penalties for surpassing or missing targets are symmetric.

To illustrate the distinction between grant-date and realized pay, suppose that a CEO's compensation in 2010 and 2011 consisted of a salary of $500,000 paid each year, and 50,000 shares of restricted stock awarded at the beginning of 2010 that become non-forfeitable ("vest") at the end of 2011. Suppose further that the company's stock price rose from $10 to $30 over the course of these two years. This CEO's *grate-date* pay (which includes the grant-date value of the restricted stock) was $1000,000 in 2010 (consisting of the 2010 $500,000 base salary and the unvested stock with a grant-date value of $500,000) and the 2011 salary of $500,000. But, his *realized* pay (consisting of his base salary plus the amount realized upon vesting) was $500,000 in 2010 and $2000,000 in 2011 ($500,000 in base salary plus $1500,000 from the vesting of his stock at the end of 2011).

Grant-date and realized pay are both legitimate measures of CEO compensation and each is a legitimate answer to a different question. Compensation committees evaluating the competitiveness of their CEO pay package at the beginning of the year (that is, before performance results are tallied) should focus on grant-date pay levels. In contrast, realized pay levels will (by definition) depend on the company's current and past performance, and are therefore most useful in evaluating whether ultimate rewards have been commensurate with company performance.

The distinction between grant-date and realized pay is also critical for researchers estimating the link between pay and performance. For example, researchers beginning with (I confess, reluctantly) Murphy (1985) have assessed the relation between pay and performance by regressing total grant-date compensation on measures of corporate performance (using CEO fixed-effects or first-differences to control for unobservable factors affecting pay levels). However, consider two otherwise identical executives, the first paid $1 million annually in base salary and the second paid $1 million annually in restricted shares. Researchers regressing grant-date pay levels on performance would conclude that neither executive is paid for performance, when in fact the second CEOs realized pay is strongly related to performance.

The SEC has helped confuse the distinction between grant-date and realized compensation by conflating elements of each in the "Summary Compensation Table" required in corporate proxy statements. In particular, since 2009, the SEC has required

companies to report the grant–date fair-market values of stock and option grants in the Summary Compensation Table, while at the same time reporting the realized (rather than target) payouts from non–equity-based bonus plans. In addition, the SEC rules are particularly confusing for companies that pay annual bonuses partly in cash and partly in stock and options, as is common in financial services. As an example, suppose that a CEO receives a bonus of $10 million in January 2012 for performance in 2011, and that $4 million is paid in cash and the remaining $6 million in stock and options. According to SEC rules, the $4 million cash bonus is included as part of 2011 compensation (and reported in the firm's 2012 proxy statement), while the $6 million bonus paid in the form of stock and options is included as part of 2012 compensation (and not reported until the firm's 2013 proxy statement).

Adding to the confusion between grant–date and realized pay was the (thankfully temporary) existence of a *third* measure mandated by the SEC and included in the Summary Compensation Table in proxy statements issued between 2007 and 2009 (covering compensation paid between 2006 and 2008). Under the SEC's 2007–2009 reporting requirements, "SEC Total Compensation" included the *accounting expense* the company records for stock and options during the year under Financial Accounting Standard 123R (FAS 123R) discussed below in Section 3.8.4. Using the accounting expense for valuing options instead of the grant–date value of options was a last–minute change to the reporting requirements made by the SEC in December 2006 without public comment. Under the SEC approach that mandates the use of accounting numbers in the table, the grant–date value of the $500,000 grant vesting in two years is reported as $250,000 in the grant year and $250,000 in the following year—numbers that bear no meaningful economic relationship to anything in the system. Fortunately, the confusion was relatively short-lived: in late 2009 the SEC revised its disclosure rules to include grant–date values rather than annual accounting expenses in the summary pay table.

Another element of the confusion in describing the typical CEO pay package reflects the statistical distinction between averages and medians. Suppose, for example, that there are eleven CEOs in an industry, ten receiving compensation of $1 million and the eleventh receiving $12 million. The *average* compensation in this industry is $2 million (calculated by summing all compensation amounts and dividing by 11), while the *median* is only $1 million (calculated as the compensation where half the CEOs are paid more and half the CEOs are paid less). Average and median pay are, again, both legitimate measures of CEO pay, but are answers to different questions. Average pay is relevant in assessing aggregate levels of pay (a reader can multiply the average pay by the number of CEOs and get total compensation paid to all CEOs), while median pay is more relevant in describing compensation for a "typical" CEO.

Figure 1 illustrates the 2011 grant–date and realized compensation for CEOs in firms listed in Standard and Poors S&P 500 (essentially the largest 500 US firms ranked by market value). The data are based on proxy statement information reported in

Standard & Poors' ExecuComp database for the 465 S&P 500 firms.[4] For both measures, total compensation is comprised of six basic components: (1) base salaries; (2) discretionary bonuses; (3) non-equity incentives (based on both annual and multi-year performance measures); (4) stock options; (5) stock awards; and (6) other pay.[5] Base salaries and the payouts from discretionary (non-formulaic) bonuses are the same for both grant-date and realized total compensation. However, the definitions of the remaining pay components vary with the measure utilized.

For *grant-date pay*, non-equity incentives are evaluated at the target level of payout (or, calculated as the average of the minimum and maximum payout if the target is not reported).[6] The grant-date value of stock options is defined as the company's estimate of the present value of the options on the grant-date: this value is typically based on Black and Scholes (1973) or similar methodologies and approximates the amount an outside investor would pay for the option. Similarly, the grant-date value of stock awards is calculated as of the grant date using the grant-date market price, which in turn approximates the amount an outside investor would pay for the stock. "Other compensation" includes perquisites, signing bonuses, termination payments, and above-market interest paid on deferred compensation. In addition, "other compensation" includes the change in the actuarial value of pension benefits, which typically constitutes a large percentage of compensation for those executives with supplementary defined-benefit pension plans.[7]

For *realized pay*, non-equity incentives are defined as actual payouts during the fiscal year, including both amounts paid in formula-based annual bonus plans, and current-year payouts from longer-term plans. Stock options are calculated as the gains realized by exercising options during the year, and stock awards are calculated as the value of

---

[4] I adopt the convention that companies with fiscal closings after May 31, in year "T" are assigned to fiscal year "T" while companies with fiscal closings on or before May 31, Year "T" are assigned to fiscal year "T−1". Thus, the 2011 fiscal year includes companies with fiscal closings between June 1, 2011 and May 31, 2012. The data in Figure 2.1 are based on the ExecuComp's May 2012 update, and exclude 35 companies that had not yet filed proxy statements by May 2012.

[5] The categories in Figure 2.1 are designed to correspond to the SEC disclosure requirements effective as of December 2006. Under the prior disclosure requirements, firms separately reported "annual bonuses" and "payouts from long-term performance plans". Under the 2006 requirements, both annual cash bonuses from short-term incentive plans and long-term performance bonuses are considered "non-equity incentive compensation" if they are based on pre-established and communicated performance targets. If they are not based on pre-established and communicated targets the SEC (and I) treat them as discretionary bonuses.

[6] The actual payouts during the year are used as an estimate for grant-date non-equity incentives in firms without reported targets or caps.

[7] The "change in the actuarial value of pension benefits" is the year-to-year change in the actuarial present value of the CEO's accumulated benefit under all defined benefit and actuarial pension plans, assuming a normal retirement age as defined in each company's plan (or, if not so defined, the earliest time at which the CEO may retire under the plan without any benefit reduction due to age). The pension information in Figure 2.1 was first available in 2006, and these amounts are therefore excluded in my historical analyses below.

**Figure 1** 2011 pay for CEOs in S&P 500 companies. *Note:*Figure 1 is based on proxy statement information compiled in Standard & Poors' ExecuComp database for 465 S&P 500 firms with fiscal closings between June 2011 and May 2012, based on ExecuComp's May 2012 update. *Grant-date Pay:* Base Salary and Discretionary Bonus reflects amounts *actually received for the fiscal year.* Non-Equity Incentives *evaluated at target level* (or average of minimum and maximum if target not reported). Stock Options *evaluated at grant-date using firm-estimated present value* (typically Black and Scholes (1973) calculations). Stock Awards *evaluated at grant-date using firm-estimated present value* (typically grant-date market price), including both time-lapse restricted stock and performance shares. Other Compensation includes perquisites, signing bonuses, termination payments, above-market interest paid on deferred compensation, and the change in the actuarial value of pension benefits. *Realized Pay:* Base Salary and Discretionary Bonus reflects amounts *actually received for the fiscal year.* Non-Equity Incentives defined as *payouts during the fiscal year* (including payouts on awards made in prior years). Stock Options defined as *gains executive realized by exercising options during the fiscal year.* Stock Awards defined as *value of awards vesting during the fiscal year* (valued on the date of vesting). Other Compensation includes perquisites, signing bonuses, termination payments, above-market interest paid on deferred compensation, and pension benefits paid during the year. The pay-composition percentages for Average Compensation are calculated as the average ratio of each component to total compensation for each CEO. The composition percentages for Median Compensation are calculated as the median ratio of each component: median ratios do not sum to 100% (because the sum of the medians is not the median of the sum).

the stock (or other equity instruments) as of the vesting date. Other compensation includes perquisites, signing bonuses, termination payments, above-market interest paid on deferred compensation, and the actual payments made to the CEO during the year under pension or retirement plans.

The first two columns in Figure 1 depict *average* grant-date and realized compensation. The pay-composition percentages are constructed by first calculating the

composition percentages for each CEO, and then averaging across CEOs. The average grant-date CEO Pay in S&P 500 firms in 2011 was $11.6 million, compared to average realized pay of $12.3 million. Stock awards are the largest single component of both grant-date and realized pay in 2011. The "Other Pay" component of grant-date pay is large compared to the corresponding component for realized pay, reflecting that the definition of grant-date pay includes the (generally positive) change in the actuarial present value of pension benefits during the year. In contrast, the realized pay for pensions include only pension benefits paid during the year for proxy-named executives (which excludes amounts to be paid after retirement).

The remaining two columns in Figure 1 depict *median* compensation. The composition percentages for median pay are calculated as the median ratio of each component: median ratios do not sum to 100% (because the sum of the medians is not the median of the sum). Median compensation is typically lower than average pay, since a small number of very-highly paid CEOs will increase the average pay but not the median pay. For example, ConocoPhillips's CEO James Mulva realized $141 million through exercising stock options in 2011. If the options had not been exercised, his pay would have fallen to "only" $5.3 million, and the average realized compensation for the 465 executives in Figure 2.1 would fall $303,000 from $12.436 million to approximately $12.133 million. Equity awards for the median executive are dominated by stock (rather than option) awards, and together option and stock awards comprise about half of total compensation for the typical executive.

The difference between grant-date and realized values, and averages and medians, is especially pronounced for stock options. Figure 2 shows the average and median grant-date values and exercise gains (i.e. realized values) for stock options granted to or exercised by CEOs in S&P 500 firms from 1992 to 2011. As shown in the figure, the average grant-date values (dotted line) and exercise gains (solid line) were remarkably similar leading up to the 2000 burst in the Internet bubble. In contrast, average exercise gains increased while average grant-date values fell leading up to the 2008 financial crisis, reflecting the shift in grants primarily reflecting the shift from options to restricted stock described in more detail below.

Figure 2 shows that median grant values and exercise gains were always below their respective averages. Interestingly, the median exercise gain was zero except for in the 2004–2007 period, indicating that less than half of the S&P 500 CEOs exercised options during most years in the sample (including 2000, when the *average* gain across all S&P 500 CEOs exceeded $12 million).[8]

---

[8] The "spike" in exercise gains in 2006 likely reflects companies accelerating the exercisability of options in anticipation of new accounting rules that would require an accounting expense for outstanding non-exercisable options; see Choudhary, Rajgopal, and Venkatachalam (2009) and the discussion in Section 3.8.4 below.

**Figure 2** Average and Median Stock Option Grant-Date Values and Exercise Gains for CEOs in S&P 500 Firms, 1992–2011. *Note:* Grant-date values based on company fair-market valuations, when available, and otherwise based on ExecuComp's modified Black–Scholes approach. Dollar amounts are converted to 2011-constant dollars using the Consumer Price Index.

Figure 3 shows how average grant-date pay for CEOs has evolved from 1970 to 2011. The data are adjusted for inflation and are based on information extracted from annual *Forbes* surveys (1970–1991) and Standard & Poors ExecuComp Database (1992–2011).[9] Non-equity pay includes base salaries, payouts from short-term and long-term bonus plans, deferred compensation, and benefits. Total compensation includes non-equity compensation plus equity-based compensation, including the grant-date values of stock options and restricted stock.[10] Due to changing reporting requirements and data availability some of the estimates of grant-date compensation are approximations,

---

[9] The *Forbes* survey includes data from the largest 500 firms ranked by market capitalization, assets, sales, and net income; the union of these sets includes approximately 800 CEOs per year. The ExecuComp survey includes data from firms in the S&P 500, S&P MidCap 400, S&P SmallCap 600, plus additional firms not in these indices, and covers approximately 1800 CEOs per year. Compustat historical data were used to identify firms included in the S&P 500 at the end of each fiscal year.

[10] ExecuComp's modifications for 1992–2006 include using 70% of the option full term, and Winsorizing dividends and volatilities. Equity compensation prior to 1978 estimated based on option compensation in 73 large manufacturing firms (based on Murphy (1985)), equity compensation from 1979 through 1991 estimated as amounts *realized* from exercising stock options during the year, rather than grant-date values. Using the amounts realized from the exercise of options (rather than the value of options granted) from 1978 to 1991 is also not expected to impose a large bias in the general trend in options and compensation. Indeed, Frydman and Saks (2005) show that trends based on grants and exercises were nearly indistinguishable during this period. In addition, Hall and Liebman (1998) analyze trends in grant-date option values during the 1980s and document a very similar pattern to that shown in Figure 3.

**Figure 3** Average Equity and Non-equity Grant-Date Pay for CEOs in S&P 500 Firms, 1970–2011. *Note:* Compensation data are based on all CEOs included in the S&P 500, using data from *Forbes* and ExecuComp. CEO total pay includes cash pay, restricted stock, payouts from long-term pay programs and the value of stock options granted (using company fair-market valuations, when available, and otherwise using ExecuComp's modified Black–Scholes approach). Average (median) equity compensation prior to 1978 estimated based on option compensation in 73 large manufacturing firms (based on Murphy (1985)), equity compensation from 1979 through 1991 estimated as amounts *realized* from exercising stock options during the year, rather than grant-date values. Non-equity incentive pay is based on actual payouts rather than targets, since target payouts were not available prior to 2006. Dollar amounts are converted to 2011-constant dollars using the Consumer Price Index.

but the trends depicted in Figure 3 are nonetheless historically representative. As shown in the figure, average grant–date compensation increased from about $1.1 million in 1970 to $10.9 million in 2011, down from a peak of $18.2 million in 2000.[11] Finally, the figure shows that most of the growth in CEO pay since 1990 is explained by the growth in equity-based pay. Indeed, stock and options constituted only a trivial percentage of pay in the early 1970s, and grew to be the dominant form of pay by the late 1990s.

---

[11] The 2010 average pay in Figure 3 ($10.9 million) is slightly smaller than the $11.6 million average in Figure 1. This difference largely reflects the fact that Figure 1 includes the change in the actuarial value of pension benefits, a component of compensation that was not disclosed or reported before 2006. Another difference—but relatively immaterial—is that Figure 1 includes the "target" rather than realized payouts from bonuses and other non–equity incentive plans; these data also became available after the 2006 revisions in disclosure rules. To maintain comparability in the time-series, Figure 3 excludes pensions and uses payouts rather than targets for bonus plans.

Figure 4 shows how both the composition and level of grant–date pay evolved from 1992 to 2010. Because of the skewness in the pay distribution (where a small number of CEOs receive unusually high levels of compensation), the median pay in Figure 4 is significantly lower than the average pay in Figure 3 in each year. The pay-composition percentages in the figure are constructed by first calculating the composition per–centages for each CEO, and then averaging across CEOs. As evident from the figure, underlying the growth in pay for CEOs since the 1990s is an escalation in stock–option compensation from 1993 to 2001 coupled with a dramatic shift away from options towards restricted stock from 2002 to 2011. In 1992, base salaries accounted for 41% of the $2.9 million median CEO pay package, while stock options (valued at grant–date) accounted for 25%. By 2001, base salaries accounted for only 18% of the median $9.2 million pay, while options accounted for more than half of pay. By 2011, options fell to only 21% of pay, as many firms switched from granting options to granting restricted stock (which swelled to 36% of pay).

In interpreting the time-series in Figure 4, it is important to recognize the selection bias inherent in the S&P 500. In particular, the firms in the index are selected by a com–mittee based primarily on market capitalization and industry representation. For example, during the 1990s the S&P 500 increased its representation of "new economy" firms, as these firms became more highly valued and a more important component of the



**Figure 4** Median Grant-date Compensation for CEOs in S&P 500 Firms, 1992–2011. *Note:* Compensation data are based on all CEOs included in the S&P 500, using data from ExecuComp. CEO grant-date pay includes cash pay, payouts from long-term pay programs, and the grant-date value of stock and option awards (using company fair-market valuations, when available, and otherwise using ExecuComp's modified Black–Scholes approach). Monetary amounts are converted to 2011-constant US dollars using the Consumer Price Index.

**Figure 5** Median Grant-date Compensation for CEOs in Firms Included in the 1992 S&P 500 Note: Compensation data are based on all CEOs included in the 1992 S&P 500, using data from ExecuComp. The sample size varies from 472 in 1992 to 260 in 2011. CEO grant-date pay includes cash pay, payouts from long-term pay programs, and the grant-date value of stock and option awards (using company fair-market valuations, when available, and otherwise using ExecuComp's modified Black–Scholes approach). Monetary amounts are converted to 2011-constant US dollars using the Consumer Price Index.

economy.[12] Indeed, the fraction of the S&P 500 comprised of new economy firms grew from 5.5% in 1992 to over 12% in 2001 (and remained at about 11% for the rest of the sample period). Since new economy firms have traditionally relied on stock options as a major component of pay (see Murphy, 2003), the increase in both the level of pay and the use of options in Figure 4 in part reflects changes in the composition of the S&P 500.

Figure 5 replicates Figure 4 after restricting the sample to only firms included in the S&P 500 in 1992. This sample restriction attenuates the increase in pay levels, which increased by 165% from 1992 to 2000 (instead of 220% as in Figure 4). The figure also suggests that CEO pay continued to increase until 2007 (a starkly different pattern than suggested by Figure 4). However, while Figure 5 mitigates the S&P 500 selection bias in Figure 4, it is subject to a survivor bias: only half of the S&P 500 firms in 1992 were still publicly traded in 2011.

---

[12] I define new economy firms as companies with primary SIC designations of 3570 (Computer and Office Equipment), 3571 (Electronic Computers), 3572 (Computer Storage Devices), 3576 (Computer Communication Equipment), 3577 (Computer Peripheral Equipment), 3661 (Telephone & Telegraph Apparatus), 3674 (Semiconductor and Related Devices), 4812 (Wireless Telecommunication), 4813 (Telecommunication), 5045 (Computers and Software Wholesalers), 5961 (Electronic Mail–Order Houses), 7370 (Computer Programming, Data Processing), 7371 (Computer Programming Service), 7372 (Prepackaged Software), and 7373 (Computer Integrated Systems Design).

While the analysis in this chapter will generally focus on S&P 500 companies, Figure 6 shows the evolution of the level and compensation for CEO pay below the S&P 500. The data, extracted from ExecuComp, include firms in the S&P MidCap 400, S&P SmallCap 600, and a small number of other firms tracked by S&P. As evident by comparing Figures 4 and 6, the level of CEO pay below the S&P 500 is considerably smaller than pay levels for S&P 500 CEOs. In addition, while median pay for S&P 500 CEOs has more than tripled from 1970–2010, pay for CEOs below the S&P 500 merely doubled. Similar to their S&P 500 counterparts, restricted stock has replaced stock options as the primary form of equity–based compensation.

### 2.1.2 The "Cost" vs. The "Value" of Incentive Compensation

In constructing measures of total compensation, it is important to distinguish between two often confused but fundamentally different valuation concepts: the cost to the company of granting the compensation and the value to an executive of receiving the compensation. Consider, for example, a company that decides to give a share of restricted stock to its CEO vesting in five years (that is, the CEO is restricted from selling the share of stock for five years, and receives the accumulated dividends [plus



**Figure 6** Median Grant-date Compensation for CEOs in non-S&P 500 Firms, 1992–2011. *Note:* Compensation data are based on all CEOs included in the S&P MidCap 400, SmallCap 600, and a small number of other non-S&P 500 firms tracked by S&P and included in the ExecuComp database. CEO grant-date pay includes cash pay, payouts from long-term pay programs, and the grant-date value of stock and option awards (using company fair-market valuations, when available, and otherwise using ExecuComp's modified Black–Scholes approach). Monetary amounts are converted to 2011-constant US dollars using the Consumer Price Index.

interest] upon vesting). Suppose further that the market price of a share of stock is $10. The economic or opportunity cost of the stock grant to the company is the amount the company could have received if it were to sell an unrestricted share to an outside investor rather than giving the restricted share to the CEO. Ignoring the probability of forfeiture and the slight dilution discount associated with issuing a new share, the company could raise $10 by selling the share to an outside investor. Thus, the company's cost of granting the share is the price of the share on the open market.

Alternatively (but equivalently), by granting the restricted share to the CEO, the company is effectively promising to deliver one share of stock to the CEO in five years. If the company had no shares available to issue, it could satisfy this contract by purchasing a share on the open market in five years at a price that might be higher or lower than $10. If the company wanted to perfectly hedge the "price risk" of its future obligation, it could purchase a share of stock in the open market today (for $10) and deliver it to the CEO in five years. Thus, again, the company's cost of granting the share is simply the price of the share on the open market.

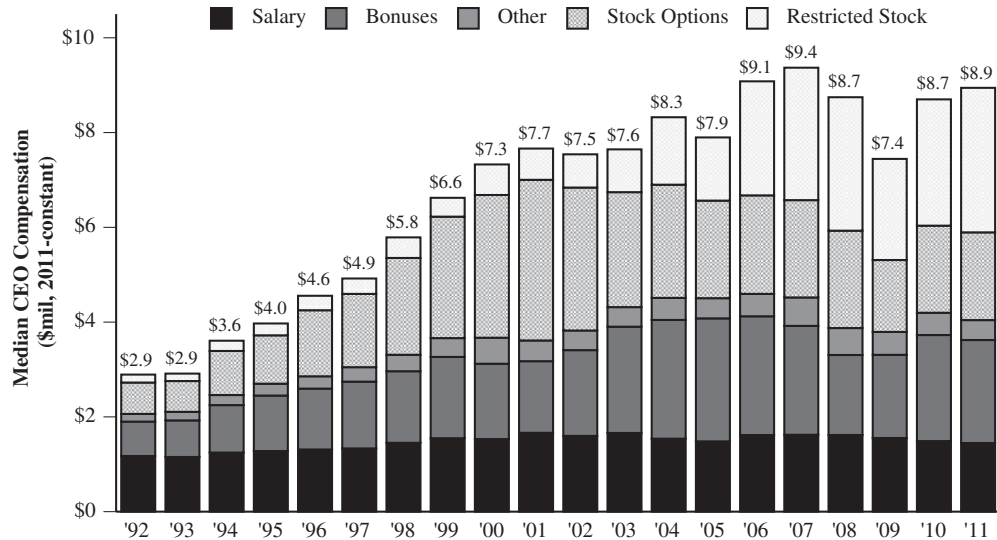But, what about the CEO? The CEO would clearly prefer to have $10 today than a promise to receive one share of stock in five years; after all, he could always take the $10 and *buy* a share of stock today, but will likely have other more-preferred uses for the $10. Moreover, if the CEO is risk averse and undiversified (in the sense that his overall wealth is positively correlated with company stock prices, through existing stock ownership, option holdings, and the risk of being fired for poor performance), the value the CEO places on the share of restricted stock will be strictly less than the fair market value of the share. Note that the CEO's value will predictably decrease as the CEO becomes more risk averse or less diversified.

Similarly, suppose that the company decides to give the CEO an option to buy a share of stock at a predetermined exercise price. The opportunity cost of granting the option is the amount an outside investor would pay for it. The outside investor is generally free to trade the option, and can also take actions to hedge away the risk of the option (such as short-selling the underlying stock). Black and Scholes (1973) and Merton (1973) demonstrated that, since investors can hedge, options can be valued as if investors were risk neutral and all assets appreciate at the risk-free rate. This risk-neutrality assumption forms the basis of option pricing theory and is central to all option pricing models, including binomial models, arbitrage pricing models, and Monte Carlo methodologies. Ignoring dilution, forfeiture, and early exercise, these now-standard methodologies provide reasonable estimates of what an outside investor would pay, and therefore measure the company's cost of granting options.

Measures of opportunity cost that ignore dilution, forfeiture, and early exercise will systematically overstate the company's cost of granting options. Dilution reduces the cost of granting options because companies typically issue new shares when options (technically, warrants) are exercised. While the impact of dilution on any specific option

grant is typically immaterial, the impact can be significant when added across all employees receiving options. Forfeiture reduces the cost because executives typically forfeit some or all of their unexercisable options upon resignation or termination.[13] Most importantly, allowing executives to exercise options before they expire reduces the company's cost of granting options because risk-averse employees—seeking diversification and liquidity—predictably exercise non-tradable options sooner than would an outside investor holding a tradable option.

However, even after appropriate adjustments for dilution, forfeiture, and early exercise, Black–Scholes values do not measure the value of the non-tradable option to a risk-averse executive. In contrast to outside investors, company executives cannot trade or sell their options, and are also forbidden from hedging the risks by short-selling company stock. In addition, while outside investors tend to be well-diversified (holding small amounts of stock in a large number of companies), company executives are inherently undiversified, with their physical as well as human capital invested disproportionately in their company. For these reasons, company executives will generally place a much lower value on company stock options than would outside investors.

Lambert, Larcker, and Verrecchia (1991) and Hall and Murphy (2002) propose measuring the value of a non-tradable option to an undiversified risk-averse executive as the amount of riskless cash compensation the executive would exchange for the option.[14] Suppose that an executive has non-firm-related wealth of $w$, holds a portfolio $S(\cdot)$ of company shares and options, and is granted $n$ options to buy $n$ shares of stock at exercise price $X$ in $T$ years. Assuming that $w$ is invested at the risk-free rate, $r_f$, and that the realized stock price at $T$ is $P_T$, the executive's wealth at time $T$ is given by

$$W_T \equiv w(l + r_f)^T + S(P_T) + n \cdot \max(0, P_T - X). \tag{1}$$

If, instead of the option, he were awarded $V$ in cash that he invested at the risk-free rate, his wealth at time $T$ would be:[15]

$$W_T^V \equiv (w + V)(l + r_f)^T + S(P_T). \tag{2}$$

---

[13] Employment agreements often provide for accelerated vesting in situations where the executive is terminated by the company without cause.

[14] Meulbroek (2001) measures the value:cost "inefficiency" of options using a completely different (non-utility-based) but complementary approach. Her method enables her to make precise estimates of what she calls the "deadweight cost" of option grants without knowledge of the specific utility function or wealth holdings of executives. Her approach produces a lower bound estimate of the value-cost inefficiency since her goal is to isolate the deadweight cost owing to sub-optimal diversification, while abstracting from any additional deadweight cost from the specific structure of the compensation contract.

[15] Cai and Vijh (2005) adopt a more-realistic (but computationally more difficult) assumption that the executive's safe wealth is optimally allocated between a riskless asset and the market portfolio. An advantage of the Cai–Vijh approach is that the certainty-equivalent values of options can never exceed Black–Scholes values.

Assuming that the executive's utility over wealth is $U(W)$, we can define the executive's value of $n$ options as the "certainty equivalent" $V$ that equates expected utilities (1) and (2):

$$\int U(W_T^V)f(P_T)\mathrm{d}P_T \equiv \int U(W_T)f(P_T)\mathrm{d}P_T. \tag{3}$$

Solving (3) numerically requires assumptions about the form of the utility function, $U(W)$, and the distribution of future stock prices, $f(P_T)$. I follow Hall and Murphy (2002) in assuming that the executive has constant relative risk aversion $\rho$, so that $U(W) \equiv ln(W)$ when $\rho=1$, and $U(W) \equiv \frac{1}{1-\rho}W^{1-\rho}$ when $\rho \neq 1$. I adopt the Capital Asset Pricing Model (CAPM) and assume that the distribution of stock prices in $T$ years is lognormal with volatility $\sigma$ and expected value equal to $(r_f + \beta(r_m - r_f) - \sigma^2/2)T$, where $\beta$ is the firm's systematic risk and $r_m$ is the return on the market portfolio.[16]

Calculating certainty equivalents from (3) requires data on stock and option grants and holdings (available from corporate proxy statements[17]), and also requires unobservable data on executive "safe wealth" (i.e. wealth not correlated with company stock prices) and executive risk aversion. Following Hall and Murphy (2002), I assume that CEOs have relative risk-aversion parameters of 2 or 3, and that each CEO has "safe wealth" equal to the greater of $5 million (in 2011-constant dollars) or four times the current cash compensation.[18] For other inputs, I assume a market risk premium of 6.5%, set the risk-free rate to the yield on 7-year US Treasuries, estimate dividend yields as the average yield over the past 36 months and volatilities based using the last 48 months of stock returns. Dividend yields above 5% are set to 5%, while volatilities below 20% or above 60% are set to 20% and 60%, respectively. As a simplifying assumption, I assume that the term for all options and restricted stock grants equals the term on the largest option grant (or five years if no options are granted), and assume that the executive's prior holdings of stock and options are fixed throughout the term of the new grants. Finally, I assume (somewhat arbitrarily) that the risk-adjusted value of accounting-based bonuses is worth 90% of target bonuses.

---

[16] For tractability, I assume that the distribution of future stock prices is the same whether the executive receives options or cash. If the grant provides incentives that shift the distribution, and if the shift is not already incorporated into stock prices as of the grant date, I will underestimate both the cost and value of the option.

[17] Under pre-2006 disclosure rules, companies reported only the aggregate number of options outstanding at the end of each year, and the intrinsic value of the in-the-money options. Following the procedure described in Murphy (1999) and adopted by Core and Guay (2002), I subtract the current-year grant from the year-end option holdings and calculate the number and average exercise price of prior grants.

[18] The results are generally robust to reasonable changes in these assumptions. In addition, for post-2006 data, I re-estimated certainty equivalents after including the actuarial value of pension benefits as safe wealth; the results are generally unaffected by this change.

Figure 7 shows the 1992–2011 evolution of risk–adjusted pay for CEOs in S&P 500 firms, assuming constant relative risk aversion of 2 or 3. The bar height depicts median pay without risk adjustments from Figure 4. Several features of the figure are worth noting:

- The value of compensation from the perspective of risk–averse undiversified CEOs can be substantially less than the cost of compensation reported in company proxy statements. For example, in 2001 (at the peak of the use of stock options), the median risk–adjusted pay for CEOs with constant relative risk aversion of 3 ($2.6 million) was less than one third of the median reported pay ($9.3 million).
- While reported pay levels increased significantly between 1998 and 2001 (driven primarily by the escalation in the grant-date values of stock options), risk–adjusted pay actually *fell* over this time period as a larger percentage of pay was being delivered in the form of risky stock options.
- Similarly, while reported pay levels were relatively flat from 2002 to 2007, risk–adjusted pay grew substantially as risky stock options were increasingly replaced by less–risky stock awards.



**Figure 7** Median Risk-Adjusted Pay for CEOs in non-S&P 500 Firms, 1992–2011. *Note:* Risk-adjusted pay is estimated using the "certainty equivalence" approach, estimated numerically assuming that the executive has constant relative risk aversion (rra) of 2 or 3, and assuming (using the Capital Asset Pricing Model) that the distribution of stock prices over the actual term of the options granted is lognormal with volatility $\sigma$ and expected value $(r_f + \beta(r_m - r_f) - \sigma^2/2)T$, where $\sigma$ and $\beta$ are determined using monthly stock-return data over 48 months, $r_f$ is the country-specific average yield on government securities during the year of grant, and $r_f - r_m = 6.5\%$ is the market risk premium. assuming relative risk aversion of 2 or 3; safe wealth is assumed to be the greater of $5 million or four times total compensation (in 2011-constant dollars). The risk-adjusted value of accounting-based bonuses is assumed to be worth 90% of actual bonuses.

The qualitative results in Figure 7 are robust to alternative definitions of risk aversion, safe wealth, equity premiums, and option terms. Calculating more precise estimates of risk-adjusted compensation for individual CEOs requires unavailable data on outside wealth and unobservable measures of individual risk aversion. In addition, more-precise estimates should allow CEOs to invest outside wealth in the market portfolio (Cai and Vijh, 2005) and allow for early exercise and different vesting and exercise terms of current grants and existing holdings. Nonetheless, the results in Figure 7 highlight that inferences based on reported grant-date compensation do not necessarily extend to risk-adjusted compensation.

## 2.2  Measuring Executive Incentives

Conceptually, the incentives created by any compensation plan are determined by two factors: (1) how performance is measured; and (2) how compensation (or wealth) varies with measured performance. Most of the executive compensation literature has focused on the relation between CEO and shareholder wealth (or, what Jensen and Murphy (1990b) defined as the "pay-performance sensitivity"), where CEOs with higher pay-performance sensitivities are defined as having better incentives to create shareholder value. Therefore, I begin this section with an analysis of different ways to measure the incentives that executives have to increase shareholder wealth. Next, given the recent focus on excessive risk-taking which many believe contributed to the financial crisis, I consider two measures of the incentives that executives have to increase stock-price volatilities. Finally, I discuss a variety of other incentive problems not neatly encapsulated in pay-performance or pay-volatility sensitivities, such as incentives to smooth or manage earnings or to pursue short-run profits at the expense of long-run value.

### 2.2.1  The Relation Between CEO and Shareholder Wealth

Most research on CEO incentives has been firmly (if not always explicitly) rooted in agency theory: compensation plans are designed to align the interests of risk-averse self-interested CEOs with those of shareholders. Following this framework, most of the focus has been on the relation between CEO compensation (or CEO wealth) and changes in firm value. Researchers have often used the ratio of equity-based total compensation to total compensation as a measure of incentives. However, the most direct linkage between CEO and shareholder wealth comes from the CEO's holdings of stock, restricted stock, and stock options. CEO wealth is also indirectly tied to stock-price performance through accounting-based bonuses (reflecting the correlation between accounting returns and stock-price performance), through year-to-year adjustments in salary levels, target bonuses, and option and restricted stock grant sizes, and through the threat of being fired for poor stock-price performance. The CEO pay literature has yet to reach a consensus on the appropriate methodologies and metrics to use in evaluating the "indirect" relation between CEO pay and company stock-price performance.

For practical purposes, however, Hall and Liebman (1998) and Murphy (1999) show that virtually all of the sensitivity of pay to corporate performance for the typical CEO is attributable to the direct rather than the indirect part of the CEO's contract, and the direct part can be measured from information available in corporate proxy statements.

Since agency costs arise when agents receive less than 100% of the value of output, the CEO's share of ownership is a natural measure of the potential severity of the agency problem. In particular, the CEO's percentage holdings of his company's stock measures how much the CEO gains from a $1 increase in the value of the firm, and how much he loses from a $1 decrease. Computing percentage ownership for restricted and unrestricted shares is trivial (simply divide by the total number of shares outstanding). Including stock options in a percentage holdings measure is more complicated, since options that are well out-of-the-money provide few incentives to increase stock prices, while options that are well in-the-money provide essentially the same incentives as holding stock. Therefore, each stock option should count somewhat less than one share of stock when adding the holdings to form an aggregate measure of CEO incentives, and the "weight" should vary with how much the option is in (or out) of the money. In constructing an aggregate measure of CEO incentives, I weight each option by the "Option Delta", defined as the change in the value of a stock option for an incremental change in the stock price. Option Deltas range from near zero (for deep out-of-the-money options) to near one (for deep in-the-money options on non-dividend paying stock).[19] I call our measure the "effective ownership percentage" to distinguish it from the actual ownership percentage based only on stock (and not option) holdings.

Figure 8 shows the evolution of the median effective percentage ownership for CEOs in S&P 500 firms from 1992 to 2011. The percentage ownership for stock and restricted stock is calculated by dividing the CEOs shareholdings by the total number of shares outstanding. Effective percentage ownership for stock options is measured by weighting each option held by the executive at the end of the fiscal year by "Option Delta" for that option (which varies according to the exercise price and time remaining

---

[19] The percentage option holdings multiplied by the option delta is a measure of the change in CEO option-related wealth corresponding to a change in shareholder wealth. More formally, suppose that the CEO holds $N$ options, and suppose that shareholder wealth increases by $1. If there are S total shares outstanding, the share price $P$ will increase by $P = \$1/S$, and the value of the CEO's options will increase by $N P(\partial V/\partial P)$, where $V$ is the Black–Scholes value of each option, and $(\partial V/\partial P)$ is the option delta. Substituting for $P$, the CEO's share of the value increase is given by $(N/S)(\partial V/\partial P)$, or the CEO's options held as a fraction of total shares outstanding multiplied by the "slope" of the Black–Scholes valuation. For examples of this approach see Jensen and Murphy (1990a), Yermack (1995), and Murphy (1999). Hall and Murphy (2002) offer a modified approach to measure the pay-for-performance incentives of risk-averse undiversified executives. An alternative approach, adopted by Jensen and Murphy (1990b), involves estimating the option pay-performance sensitivity as the coefficient from a regression of the change in option value on the change in shareholder wealth.

**Figure 8** Median Effective Percentage Ownership for CEOs in S&P 500 Firms, 1992–2011. *Note:* Percentage ownership for stock and restricted stock measured as the CEO's shareholdings divided by the total number of shares outstanding. Effective percentage ownership for stock options measured by weighting each option held by that options "Black–Scholes Delta" and dividing by the total number of shares outstanding. Year-end options under the pre-2006 disclosure rules estimated using the procedure described in Murphy (1999).

to exercise), and dividing by the total number of shares outstanding.[20] As shown in the figure, stock and restricted stock holdings for the median S&P 500 executive has grown modestly over the 20-year period (reflecting the increased popularity of restricted stock), ranging from 0.12% to 0.15%. Over the same time period, total effective ownership (including delta-weighted options) doubled from 0.35% in 1992 to 0.69% in 2003, before falling to 0.38% in 2011. The drop in ownership in 2008 depicted in Figure 8 primarily reflects that most options held by CEOs at the end of 2008 were substantially out-of-the-money and therefore had low incentives and low Option Deltas.

The measure of effective CEO ownership in Figure 8 is essentially the "Pay-Performance Sensitivity" introduced by Jensen and Murphy (1990b). The primary difference is that I am measuring the effective ownership percentage, while Jensen and Murphy measured the change in CEO wealth per $1000 change in shareholder wealth,

[20] Proxy disclosure rules effective since December 2006 provide the details on year-end option portfolios required to estimate Options Deltas. Year-end portfolios prior to 2006 are estimated using the procedure described in Murphy (1999) and adopted by Core and Guay (2002).

which equals the effective ownership percentage multiplied by ten. The other difference is that Jensen and Murphy also include indirect incentives from cash compensation and disciplinary terminations. Using data from 1974 to 1986, Jensen and Murphy estimate a median pay-performance sensitivity for stock and options of $2.50 for every $1000 change in shareholder wealth, which corresponds to an ownership percentage of 0.250%.[21] Therefore, by the end of 2003, pay-performance sensitivities had nearly tripled the data from 1974 to 1986. But, by year-end 2011 the pay-performance sensitivity was slightly above its 1992 level, or about 50% higher than the Jensen–Murphy estimate.

The average market capitalization of firms in the S&P 500 grew (in 2011-constant dollars) from $10.0 billion in 1992 to $35.8 billion in 2000 (before falling to $22.7 billion in 2011), therefore the dollar value of the typical CEOs ownership position is large even if his percentage holding is low. Hall and Liebman (1998) argue that a better way to measure CEO incentives is as the change in CEO wealth for a 1% change in the value of the firm rather than as the ownership percentage. Baker and Hall (2004) provide some theoretical justification for using this measure. In particular, Baker and Hall show that percentage ownership is the right measure of incentives when the marginal product of the CEO effort is constant across firm size, such as a CEO contemplating a new corporate headquarters that will benefit the CEO but perhaps not the shareholders, or an outside takeover bid that will benefit outside shareholders but perhaps not the CEO. But, the Hall-Liebman measure is appropriate when the marginal product of the CEO effort scales with firm size, such as a corporate reorganization (assuming it takes the same amount of CEO effort to reorganize a big firm as a small firm).

Figure 9 shows the evolution of the Hall-Liebman measure—what Frydman and Jenter (2010) call "equity at stake"—from 1992 to 2011. The equity-at-stake measure is calculated as 1% of the effective ownership percentage multiplied by the firm's market capitalization.[22] In 1992, each 1% change shareholder wealth resulted in a $181,000 change in CEO wealth for the median CEO in the S&P 500. The equity-at-stake measure grew to nearly $900,000 in 2000 and again in 2005, before plummeting to $265,000 in 2008 as a result of both the decline in market capitalizations and the decline in Option Deltas.

As an alternative to both the Jensen-Murphy and Hall-Liebman measures, Edmans, Gabaix, and Landier (2009) provide theoretic justification for measuring incentives

---

[21] Including incentives from potential dismissals and performance-related changes in the value of salaries, bonuses, and option grants, increased the "final" Jensen–Murphy estimate to $3.25 per $1,000, or an effective ownership percentage of 0.325%.

[22] Suppose that the CEO holds M shares and N options. If the share price $P$ increases by 1%. If there are $S$ total shares outstanding, the value of the CEO's portfolio will increase by $0.01P(M + N(\partial V/\partial P))$ or $0.01(PS)[(M + N(\partial V/\partial P))/S]$, where $PS$ is the firm's market capitalization and the quantity in the square brackets is the equation for the CEO's effective ownership percentage.

**Figure 9** Median Equity at Stake for CEOs in S&P 500 Firms, 1992–2011. *Note:* Following Frydman and Jenter (2010), Equity-at-Stake is measured as the effective ownership percentage multiplied by 1% of the firms market capitalization (in thousands of 2011-constant dollars).

using the "wealth–performance elasticity" (i.e., the percentage change in CEO wealth corresponding to a percentage change in firm value) when the CEO effort has a multiplicative (rather than additive) effect on both CEO utility and firm value. In practice, creating this measure generally requires data not available to researchers (in particular, the CEO's wealth beyond his portfolio of company stock and options).[23]

### 2.2.2 The Relation Between CEO Wealth and Stock-Price Volatilities

Suspicions that executive compensation policies in financial services firms contributed to the 2008-2009 financial crisis eventually broadened to similar suspicions for

---

[23] Several early studies, including Murphy (1985) and Gibbons and Murphy (1992), used the "pay-performance elasticity", defined as the percentage change in current compensation associated with a percentage change in company performance. While the pay-performance elasticity reflects how boards adjust current compensation to changes in performance, it ignores the CEO's portfolio of stock and options and therefore does not measure CEO incentives. Edmans, Gabaix, and Landier (2009) suggest a measure where the change in CEO wealth from stock and option holdings is divided by the CEO's current compensation rather than the CEO's total wealth; this measure is proportional to the wealth-performance elasticity to the extent that CEO wealth is proportional to current compensation. As emphasized by Murphy (1999), the empirical advantage of elasticity measures is that they are typically independent of firm size. In contrast, the Jensen-Murphy "effective ownership" percentage is predictably smaller for CEOs of larger firms.

companies outside the financial sector. In December 2009, as part of the continued fallout from the crisis, the SEC began requiring all publicly traded companies to disclose and discuss compensation policies and practices that might provide incentives for executives to take risks that are reasonably likely to have a materially adverse effect on the company.

When executives receive rewards for upside risk, but are not penalized for downside risk, they will naturally take greater risks than if they faced symmetric consequences in both directions. For top executives, rewarded primarily with equity-based compensation, the primary source of risk-taking incentives emanates from stock options. The pay-performance relation implicit in stock options is inherently convex, since executives receive gains when stock prices exceed the exercise price, but their losses when the price falls below the exercise price are capped at zero.

Since equity is a "call option" on a leveraged firm (Black and Scholes, 1973), equity-based pay in a leveraged firm can provide similar risk-taking incentives to those provided by stock options in an all-equity firm. Consider, for example, an investment opportunity promising equal chances of a $400 million gain and a $600 million loss (i.e. a net-present value of –$100). Shareholders in a $1 billion all-equity firm will have no incentive to pursue this negative NPV investment, because they will bear 100% of both the gains and losses. But, suppose the firm has only $100 million in equity, and $900 million in debt. Equity holders receive 100% of the upside, but their downside liability is limited to the value of their initial equity stake ($100 million). Thus, from the perspective of the equity holders, the project has a net present value of +$150 million.

The conflict of interest between shareholders and debtholders—dubbed the "Agency Cost of Debt" by Jensen and Meckling (1976)—has led several researchers to measure risk-taking incentives by leverage ratios and to prescribe CEO pay structures that include debt as well as equity.[24] However, it is worth noting that it is not leverage *per se* that creates risk-taking incentives, but rather the limited liability feature of equity. For example, the shareholders in the example in the prior paragraph would have incentives to take the negative NPV project even if the firm was a $100 million all-equity firm; in this case losses greater than $100 million would be borne by the government or society, etc., and not by debtholders. It is also worth noting that the severity of the risk-taking incentives depends on the maximum downside risk compared to the dollar amount of equity, and not the value of equity compared to the overall value of equity plus debt. The *level* of debt is important only to the extent that is available to fund risky negative NPV projects.

Since the value of a stock option (or the value of equity in a leveraged firm) increases monotonically with stock-price volatilities, options (and limited liability)

---

[24] See, for example, Sundaram and Yermack (2007), Edmans and Liu (2011); Edmans, "How to Fix Executive Compensation", Wall Street Journal (2012). "Debt compensation" typically consists of deferred compensation or nonqualified defined-benefit pension plans, where the executive joins other unsecured creditors in bankruptcy.

provide incentives for executives to increase such volatilities. In Section 2.2.1, the cal-culations for pay–performance sensitivities for stock options depended on the Option Delta, defined as the change in the value of a stock option associated with an incre-mental change in the stock price. Similarly, the calculations for pay–volatility sensitivities for stock options depend on the Option Vega, typically defined as the change in the value of a stock option associated with one percentage-point increase in the stock-price volatility (e.g. from 30% to 31%). Option Vegas are typically highest when stock prices are near the option's exercise price.

Following Fahlenbrach and Stulz (2011)'s analysis of executive compensation and the financial crisis, I consider two option-based measures for incentives to increase stock-price volatilities:

$$\text{Total Option Vega} = \text{Change in value of outstanding options for a one percentage-point increase in volatility.}$$

$$\text{Vega Elasticity} = \text{Percentage change in value of outstanding options for a one percentage-point increase in volatility.}$$



**Figure 10** Option-based incentives to increase volatility by CEOs in S&P 500 Firms, 1992–2011. *Note:* The Total Option Vega is defined as the change in value of outstanding options for a one percentage-point increase in volatility. Vega Elasticity is defined as the percentage change in value of outstanding options for a one percentage-point increase in volatility.

518

Figure 10 shows the time trends in the two measures of pay-volatility sensitivities for the median executive in a S&P 500 firm from 1992 to 2011. The left-hand axis reports the Total Option Vega, which reached its peak in 2003 (when the median CEO gained $243,000 by increasing volatility by one percent), and plummeted in 2008 to $127,000 for a one percent increase in volatility. The right-hand axis reports the percentage change in option values associated with a one percent increase in volatility. This "Vega Elasticity" remained relatively constant from 1992 to 2007 at around 1.0 (indicating that a one percentage-point increase in volatility would increase the value of CEO option holdings by about 1%). The Vega Elasticity jumped to over 5.0% in 2008, falling to 2.0% by 2011.

The differences in the two measures in Figure 10 reflect the effect of stock-market movements and, in particular, the market crash at the end of 2008 and the partial rebound by 2011. When stock prices fell (as they did abruptly in 2008, across all sectors of the economy), the options fell out of the money, which implies that the Option Vega for each option becomes smaller (remember that the Option Vega is highest when the stock price is close to the exercise price). But, it turns out that, as stock prices fall, the value of the options held fall even faster than the Option Vega. As a result, the value of options that are out-of-the-money increases more in percentage terms (but less in dollar or euro terms) as volatility increases.

One troublesome fact apparent from Figure 10 is that the two vega measures—both legitimate measures for risk-taking incentives—move in opposite directions in market downturns. There is no accepted methodology for measuring incentives for risk in executive option portfolios, or in executive equity positions in leveraged firms, or in executive contracts more generally.[25] Until the recent financial crisis—when compensation policies were blamed for contributing to the meltdown—there had been little focus on the role of compensation policies in providing incentives to take risks.

Finally, while the current controversy over executive incentives has focused on excessive risk taking, it is worth noting that the challenge historically has been in providing incentives for executives to take *enough* risk, not too much risk. Executives are typically risk-averse and undiversified with respect to their own companies' stock-price performance. On the other hand, shareholders are relatively diversified, placing smaller bets on a larger number of companies. As a result, executives will inherently be "too conservative" and want to take fewer risks than desired by shareholders. Stock options

---

[25] Although there is little theoretical guidance on the appropriate measure of risk-taking incentives, Alex Edmans (in private correspondence) suggests that the appropriate measure likely depends on the CEO's cost of increasing volatility. In particular, the Total Option Vega is likely the correct measure if the cost to increase volatility has an additive effect on CEO utility, while the Vega Elasticity is likely correct if the cost has a multiplicative effect. Dittmann and Yu (2011) propose an alternative measure related to the ratio of vega to delta.

(or other plans with convex payouts) have long been advocated as ways to mitigate the effects of executive risk aversion by giving managers incentives to adopt rather than avoid risky projects (see, for example, Hirshleifer and Suh, 1992). Similarly, there is a long history of attempts to document an empirical relation between such convexities and actual risk-taking incentives, and the results have been relatively modest.[26]

## 2.3  (Dis)Incentives from Bonus Plans[27]

Most discussions about incentives for US CEOs focus exclusively on equity-based incentives, since changes in CEO wealth due to changes in company stock prices dwarf wealth changes from any other source (Hall and Liebman, 1998; Murphy, 1999). However, from a behavioral perspective, annual and multi-year bonus plans based on accounting measures may be as important as equity in actually directing the activities of CEOs and other executives. Consider the following:

- Incentive plans are effective only if the participants understand how their actions affect the payoffs they will receive and then act on those perceptions. While CEOs likely understand how to increase accounting income (by increasing revenues and decreasing costs of goods sold), they often do not understand how their actions affect company stock prices. Therefore, bonus plans may well provide stronger incentives than equity-based plans, even though their magnitude is smaller.

- Most bonus plans are settled in cash soon after the results are tallied (e.g. after the year-end audited financials). The immediacy and tangibility of these cash awards may well provide stronger incentives than the distant and uncertain paper gains in unvested equity plans.

Unfortunately, while CEOs may indeed be motivated by their bonus opportunities, they are not necessarily motivated to increase firm value. The problems lie in the design of the typical bonus plan, illustrated in Figure 11. Under the typical plan, no bonus is paid until a lower performance threshold or hurdle is achieved, and a "hurdle bonus" is paid at this lower performance threshold. The bonus is usually capped at an upper performance threshold; after this point increased performance is not associated with an increase in the bonus. The thresholds are routinely determined by the firm's annual budgeting process. The range between the lower and upper performance thresholds (labeled the "incentive zone" in the figure), is drawn as linear but could be convex (bowl-shaped) or concave (upside-down bowl-shaped). The "pay-performance relation" (denoted by the heavy line) is the function that shows how the bonus varies throughout the entire range of possible performance outcomes.

---

[26] DeFusco, Johnson, and Zorn (1990) find some evidence that stock-price volatility increases, and traded bond prices decrease, after the approval of executive stock option plans. Similarly, Agrawal and Mandelker (1987) find some evidence that managers of firms whose return volatility is increased by an acquisition have higher option compensation than managers whose volatility declined.

[27] This section draws heavily from Murphy (1999) and Murphy and Jensen (2011).

**Figure 11** A Typical Bonus Plan Note: Under a typical bonus plan, a performance target and a target bonus for meeting that performance are set. Upper and lower performance thresholds are established which create an incentive zone within which the bonus increases with performance. Bonuses do not vary with performance outside the range established by the Lower and Upper Performance thresholds. A Hurdle Bonus is often paid when the executive reaches the lower performance threshold. The bonus can increase linearly with performance in the incentive zone (as shown here) or it can increase at a decreasing rate or an increasing rate (that is, the line can be convex or concave).

In spite of substantial variability across companies and industries, short–term and long–term bonus plans can be characterized in terms of the three basic dimensions suggested by Figure 11: performance measures, performance thresholds (that is, targets, benchmarks, or standards), and the structure of the pay–performance relation. Design flaws in any of these dimensions can provide incentives to withhold effort, to shift earnings and cash flow unproductively from one period to another (or otherwise manipulate earnings), to use capital inefficiently, and to destroy information critical to the effective coordination of disparate parts of large complex firms.

### 2.3.1  Problems with Non-Linear Pay-Performance Relations

Researchers have long acknowledged that non-linear incentive plans cause predictable problems.[28] For example, executives capable of producing well above the upper performance

---

[28] The pioneering empirical paper is Healy (1985), who found that executives use discretionary accrual charges to shift earnings to a later period whenever performance exceeds the upper performance threshold. Holmstrom and Milgrom (1987) provide the classic theoretical justification for linear contracts based on specific modeling assumptions; Edmans and Gabaix (2011) provide more general conditions for linearity.

threshold in Figure 11 have incentives to stop producing once they "max out" on their bonuses. In addition, they will do their best to transfer performance results that could have been realized this period into a later period.

Similarly, but potentially worse, is the effect of the discontinuity at the lower performance threshold in Figure 11. Executives who believe they cannot achieve at least this level of performance this year will either stop producing or "save" performance for next year by delaying revenues or accelerating expenses. Moreover, if executives see that they are not going to make the bonus pool this year, they are better off taking an even bigger hit this period (since there is no bonus penalty for missing the lower threshold by a lot instead of a little) so they can do even better next period—what accountants have called the "big bath" phenomenon. On the other hand, executives who are struggling to make the lower threshold, but still believe they can make that threshold, have incentives (provided by the threshold bonus) to do whatever is necessary to achieve the lower threshold. Their actions commonly include destroying value by loading the distribution channel so as to recognize revenues earlier, unwisely reducing R&D and required maintenance expenditures, and (in some cases) outright accounting fraud. Each of these actions shifts reported profits from next period to the current period, but does so at an unnecessary cost to the firm.

In both of these cases, the non-linearities provide incentives for CEOs to "manage earnings". In particular (and assuming that performance is measured by earnings), the bonus plan in Figure 11 provides incentives to "smooth earnings" (by shifting earnings from next period when below the lower threshold and shifting earnings to next year when above the upper threshold), while occasionally taking a "big bath" (when it is not possible, even with manipulation, to get earnings above the lower threshold).

In addition to earnings management, non-linearities also affect risk-taking behavior. In particular, when the pay–performance relation is concave (so that lower performance is penalized more than higher performance is penalized), executives can increase their total bonus payouts by reducing the variability of their performance. Conversely, convex pay–performance relations increase risk-taking incentives. Financial economists have suggested that boards purposely add convexity to CEO pay contracts to offset the reluctance of risk-averse CEOs to invest in risky (but profitable) projects.[29] More recently, some academics (as well as Congress and the popular press) have alleged that convexities in banking bonuses (where positive performance is rewarded, but negative performance is not penalized) led to excessive risk-taking that, in turn, facilitated the 2008–2009 financial crisis.

The problems with non-linearities are mitigated by eliminating caps on the upside, and finding ways to implement and enforce "negative" bonuses on the downside.[30] While it is difficult to force CEOs to write checks back to the company after a bad year, negative bonuses can be partially implemented by basing pay on multi-period cumulative

---

[29] Classic papers include Hirshleifer and Suh (1992) and Guay (1999).
[30] See Murphy and Jensen (2011) for an extended discussion and example of these practices.

performance (Holmstrom and Milgrom, 1987) or by deferring current compensation into bonus banks that can be used to fund future negative bonuses (Stewart, 1991). Another indirect way to impose negative bonuses is by reducing base salaries and offering enhanced bonus opportunities (through reduced bonus thresholds).

### 2.3.2  Problems with Performance Benchmarks

Bonuses are usually not, in practice, based strictly on a performance measure, but rather on performance measured relative to a performance benchmark (Murphy, 2000). Examples include net income measured relative to budgeted net income, EPS vs. last year's EPS, cash flow vs. a charge for capital, performance measured relative to peer-group performance, or performance measured against financial or non-financial strategic "milestones". Performance targets (one form of benchmark) typically correspond to the level of performance required to attain the executive's "target bonus".

When bonuses are based on performance relative to a benchmark, executives can increase their bonus either by increasing performance or lowering the benchmark. Performance benchmarks therefore create predictable problems whenever the participants in the bonus plan can affect the benchmark. For example, when benchmarks are based on meeting a budget, executives with bonuses tied to budgeted performance targets have strong incentives to low-ball the budget. Boards (and supervisors throughout the management hierarchy) understand these incentives and generally push for higher budgets than those suggested by executives. The result is a familiar and predictable "budget game" that ultimately destroys the information critical to coordinating the disparate activities of a large complex organization (Jensen, 2003).

As another example, when benchmarks are based on prior performance (such as bonuses based on growth or improvement), plan participants understand that increased performance this year will be penalized by higher benchmarks the next year, and will naturally take account of these dynamics when deciding how hard to work and what projects to undertake in the current year. Similarly, when bonuses are based on performance measured relative to that of colleagues, participants can increase their bonuses by sabotaging co-executives (Lazear, 1989; Gibbons and Murphy, 1990). Benchmarks based on industry peers provide incentives for selecting "weak" industries or peers, or staying too long in a defective industry (Dye, 1992).

The problems with benchmarks based on budgets, prior-year performance, co-workers, and other internally manipulable measures can be mitigated by "externalizing" the benchmark; that is, by basing the benchmark on objective measures beyond the direct control of the plan participants. In Murphy (2000), I showed that companies using external benchmarks (which I defined as benchmarks based on fixed numbers or schedules, industry performance, or the cost of capital) were less likely to manage fourth-quarter earnings than were companies with internal benchmarks. However, I was unable to explain satisfactorily cross-sectional differences in the use of internal and

external benchmarks, or why nearly 90% of the sample of 177 firms based benchmarks on budgets or prior-year performance.

### 2.3.3  Problems with Performance Measures

The problem of inappropriate performance measures is illustrated succinctly by the title of Steven Kerr's famous 1975 article, "On the folly of rewarding A, while hoping for B" (Kerr, 1975). Paying salespeople commissions based on revenues, for example, provides incentives to increase revenues regardless of the costs or relative margins of different products. Likewise, paying rank-and-file workers "piece rates" based on units produced provides incentives to maximize quantity irrespective of quality, and paying a division head based solely on divisional profit leads that division head to ignore the effects of his decisions on the profits of other divisions. Similarly, paying CEOs based on short-run accounting profits provides incentives to increase short-run profits (by, for example, cutting R&D) even if doing so reduces value in the long run.

Conceptually, the "perfect" performance measure for a CEO is the CEO's personal contribution to the value of the firm.[31] This contribution includes the effect that the CEO has on the performance of others in the organization, and also the effects that the CEO's actions this year have on performance in future periods. Unfortunately, the CEO's contribution to firm value is almost never directly measurable; the available measures will inevitably exclude ways that the CEO creates value, and include the effects of factors not due to the efforts of the CEO, or fail to reveal ways that the CEO destroys value. The challenge in designing incentive plans is to select performance measures that capture important aspects of the CEO's contributions to firm value, while recognizing that all performance measures are imperfect and create unintended side effects.

While companies use a variety of financial and non-financial performance measures in their annual CEO bonus plans, almost all companies rely on some measure of accounting profit such as net income, pre-tax income, or operating profit. Accounting profit measured over short intervals is not, however, a particularly good measure of the CEOs contribution to firm value, for several reasons. First, CEOs routinely make decisions (such as succession planning or R&D investments) that will increase long-run value but not short-run profit. Second, accounting profits (like equity-based measures)

---

[31] In his classic paper on optimal contracts, Holmstrom (1979) considers a case where the principals (i.e. the shareholders) know precisely what action they want the agent (i.e. the CEO) to take, but cannot observe whether the CEO in fact took that action. Holmstrom shows that the optimal contract will include any performance measures that are useful (or "informative") in determining whether the CEO took the prescribed action. This so-called "informativeness principle" was widely embraced by many academics who used it as the theoretical justification for analyzing performance measures used in CEO contracts. However, as emphasized in Holmstrom (1992) and implicit in Holmstrom and Milgrom (1991), the informativeness principle is not applicable in the realistic multi-tasking case where the shareholders do not know precisely what actions they want the CEO to take, and indeed entrust their money to self-interested CEOs specifically because CEOs have superior skill or information in making investment decisions.

are invariably influenced by factors outside of the control of the CEO, including the effects of business cycles, world oil prices, natural disasters, terrorist attacks, etc. Third, while the measures of accounting profits typically used in bonus plans take into account both revenues and expenses, they ignore the opportunity cost of the capital employed. The use of these accounting measures provides incentives to invest in any project that earns positive accounting profits (not just those that earn more than the cost of capital), and provides no incentives to abandon projects earning positive accounting profits that are less than those required to cover their cost of capital.

Exacerbating the problems with accounting-based performance measures in bonus plans is the fact that they are often expressed as ratios (e.g., earnings per share, return on assets, return on equity, return on capital, etc.). Executives participating in such plans can increase their bonus either by increasing the numerator (accounting profits) or by decreasing the denominator (e.g. shares, assets, equity, invested capital). For example, a CEO paid on the basis of return on capital would prefer a $100 million project earning a 40% return to a $1 billion project earning a 25% return, even though the latter creates more wealth (as long as the cost of capital is less than 22%).

## 2.4 (Dis)Incentives from Capital Markets[32]

The typical accounting-based bonus plan depicted in Figure 11 provides incentives to focus on short-run accounting returns at the expense of long-run value creation, and to manipulate or smooth earnings by unproductively shifting revenues and expenses across reporting periods. Conceptually, this problem is mitigated by shifting from accounting- to equity-based plans: if markets are efficient, then the equity markets should punish executives for playing the "earnings management game". However, equity markets can exacerbate rather than mitigate the problem, by providing executives with incentives to take actions to meet or beat analyst and market expectations for earnings or certain key performance benchmarks.

Figure 12 shows the relation between the magnitude of the quarterly abnormal stock return and quarterly earnings surprises measured by the earnings forecast error, based on 172,247 firm-quarter observations over the period 1984-2010.[33] The earnings forecast error is defined as the difference between actual announced earnings per share and the median analyst forecast for quarterly earnings thirteen trading days prior to the end of the quarter, divided by the closing stock price for the quarter. Abnormal returns reflect the cumulative return from twelve days before to one day after the earnings announcement, less the buy-and-hold return from the associated Fama-French 5x5 portolio (based on size and book-to-market ratios). Accounting data are from Compustat, returns and share prices from CRSP, and earnings forecasts are from I/B/E/S.

As shown in Figure 12, stock prices react strongly and positively to small positive earnings surprises: when a firm produces earnings that beat the consensus analyst forecast

---

[32] See Jensen and Murphy (2012) for a more-detailed treatment of the analysis in this section.
[33] The data and analysis underlying Figure 12 were generously provided by David Huelsbeck.

**Figure 12** Abnormal Stock Returns in Response to Quarterly Earnings Surprises Note: The graph plots quarterly abnormal returns for growth and value firms as a function of earnings surprise at the end of the quarter. Forecast error is measured as the earnings surprise relative to the quarter-end stock price. Data are from I/B/E/S database for the final month of the fiscal quarter for which earnings is being forecast. Each "dot" represents averages for 200 portfolios ranked by the earnings surprise. Sample size is 172,247 firm-quarter observations in the period 1984-2010. *Note:* Data and analysis provided by David Huelsbeck.

by 1% the stock price rises on average by about 5.5%. Similarly, stock prices react strongly and negatively to small negative earnings surprises: when a firm misses its forecast by 1% stock prices fall by nearly 8%. But there is not much additional stock-price reaction to larger surprises (those greater than plus or minus 1% of the stock price at the end of the final month of the fiscal quarter for which earnings are being forecast). This "S-curve" feature of stock-price responses to earning surprises has been well documented in the literature (see, for example, Skinner and Sloan, 2002; Bartov, Givoly, and Hayn, 2002).

As emphasized by Jensen and Murphy (2012), the relation between a firm's top-management team and the capital markets has resulted in an equilibrium that replicates many counterproductive aspects of budget or target-based bonus systems discussed in conjunction with Figure 11. For executives holding large quantities of stock and stock options, Figure 12 portrays the non-linear pay–performance relation that defines how meeting, beating or missing analyst forecasts affects the value of their equity-based hold-ings. In particular, executives subject to such stock price responses to quarterly earnings surprises have incentives to beat analysts forecasts by a small amount (an earnings sur-prise that amounts to no more than 1% of the quarter-end stock price), but not by *too much* because the payoff from beating the forecast by a lot is not much higher than the payoff for beating it by 1%. Note also that manipulating this quarter's earnings to miss

analyst earnings forecasts by a *lot* (e.g. by shifting revenues from this quarter to the next quarter, or moving expenses from next quarter to this quarter) also provides increased ability to executives to beat *next* quarter's earnings forecast.

Following the accounting scandals in the early 2000s, several researchers have documented that executive option and equity holdings are higher in companies that restate their earnings or are accused of accounting fraud. The results are mixed. Efendi, Srivastava, and Swanson (2007) and Burns and Kedia (2006), for example, document that firms with CEOs who have large amounts of "in-the-money" options are much more likely to be involved in restatements. Bergstresser and Philippon (2006) provide evidence that the use of discretionary accruals to manipulate reported earnings is more pronounced in firms where the CEO's potential total compensation is more closely tied to the value of stock and option holdings. Johnson, Ryan, and Tian (2009) concludes that firms accused of fraud have significantly greater incentives from unrestricted stockholdings than control firms do, and unrestricted stockholdings are their largest incentive source. Erickson, Hanlon, and Maydew (2006) find in logistic regressions that the probability of being accused of fraud by the SEC is related to stock-based compensation, but find no differences between the fraud firms and a "matched" sample of firms not accused of fraud.

Temptations to manipulate the expectations market will clearly be higher for executives holding large quantities of stock and options that can be sold or exercised before markets adjust to the "real" information. Therefore, the natural remedy to mitigate manipulation is to impose longer vesting periods on restricted stock and options and holding requirements on unrestricted stock.[34] However, there is little evidence that executives actually exercise and sell large fractions of their exercisable options or sell large fractions of their unrestricted stock holdings prior to restatements or indictments. The ominous hypothesis is that executives focused on the expectations market are not following a "pump and dump" strategy (which can be controlled by imposing longer vesting and holding requirements), but rather that they are legitimately confused about the difference between increases in the short-run stock price and true value creation.

## 3. HOW WE GOT THERE: A BRIEF HISTORY OF CEO PAY

### 3.1 Introduction

Most recent analyses of executive compensation have focused on efficient-contracting or managerial-power rationales for pay, while ignoring or downplaying the causes and consequences of disclosure requirements, tax policies, accounting rules, legislation, and the general political climate. A central theme of this study is that government intervention has been both

---

[34] See, for example, Edmans et al. (2012) and (in the context of the financial crisis) Bhagat and Bolton (2011).

a response to and a major driver of time trends in executive compensation over the past century, and that any explanation for pay that ignores political factors is critically incomplete.

As will become evident in this section, there have been two broad patterns for government intervention into CEO pay. The first pattern is aptly described as knee-jerk reactions to isolated perceived abuses in pay, leading to disproportionate responses and a host of unintended and undesirable consequences. As an example discussed below in Section 3.6.1, outrage over a single $4.1 million change-in-control payment in 1982 led to strict limitations on all golden parachutes for top executives, which in turn led to a host of unintended consequences including an explosion in the use of golden parachutes, tax gross-up provisions, and employment agreements; the rules also encouraged shorter vesting periods for stock awards and early exercise of stock options. The second pattern—best described as "populist" or "class warfare"—arises in situations where CEOs (and other top executives) are perceived to be getting richer when lower-level workers are suffering. The associated attacks on wealth in these situations gave rise to disclosure rules in the 1930s, limits on tax deductibility for CEO pay in the early 1990s, and wide-ranging pay regulations in the 2010 Dodd–Frank Act. Beyond these two broad patterns, indirect intervention in the form of accounting rules, securities laws, broad tax policies, and listing requirements have also had direct impact on the level and composition of CEO pay.

Calling this second pattern "class warfare" is a bit simplistic, since (relative to other developed economies) Americans have historically been unusually tolerant of income inequality arising from exceptional efforts, ideas, and abilities. Underlying much of the outrage—and suggestive of the managerial-power hypothesis—is the perception that executive pay is "rigged" and not reflective of productivity and not set in a competitive market for managerial services.[35] Nonetheless, it is instructive to recognize that demands to reform (or punish) CEO pay are concentrated in "third parties" angry with perceived levels of excessive pay, and not shareholders concerned about insufficient links between pay and performance.

## 3.2  Executive Compensation Before the Great Depression[36]

The history of executive compensation in the United States naturally parallels the history of executives. While the vast majority of business enterprises before 1900 were small and run by owners, a new class of "salaried middle managers" emerged in a variety of industries (such as railroads and steel) with relatively large and complex firms. However, even these larger firms were typically run by founders, descendents

---

[35] While the recent Occupy Wall Street movement is insufficiently organized to speak with a single voice, a plausible interpretation of their attack on Wall Street pay (and CEO pay, more generally) is the perception that pay is rigged; see, for example, Taibbi, "Politics: OWS's Beef: Wall Street Isn't Winning - It's Cheating", *Rolling Stone* (2011).

[36] The material in this subsection is largely drawn from Wells (2010, 2011).

of founders, or individuals with large blocks of equity: there was no obvious need for executive incentive plans that tied pay to corporate performance.

Between 1895 and 1904, nearly two thousand small manufacturing firms combined to form 157 large corporations. Management responsibility in many of these new firms shifted from owners to professional executives who had management skills but no meaningful equity stakes. Over the next two decades, the void in incentives was filled by the emergence of bonuses tied to corporate profits. By 1928, nearly two thirds of the largest industrial companies offered executive bonus plans; bonuses accounted for 42% of 1,929 total executive compensation in companies with plans (Baker, 1938). While compensation was generally modest, the highest bonuses rivaled amounts even in nominal terms not seen again until the late 1970s. For example, as discussed below, Bethlehem Steel's CEO Eugene Grace received a bonus of $1.6 million for 1929 performance (over $20 million in inflation-adjusted 2010 dollars).

In spite of the increasing magnitude of the highest CEO bonuses, executive pay was not particularly controversial during the 1920s. Part of the nonchalance reflected the fact that there were no public disclosures of pay for individual executives: the bonuses at Bethlehem Steel, for example, came to light as a result of a 1930 lawsuit unrelated to compensation. Most reports at the time were speculative, based on vague descriptions of company-wide bonus formulas that would allow estimates of aggregate but not individual bonuses. Moreover, the economy was robust, unemployment was low, and shareholder returns were high, factors that would provide a safe harbor for high executive pay for the next 90 years.

In July 1930, during a lawsuit attempting to block Bethlehem's takeover of Youngstown Sheet & Tube Co., Bethlehem Steel's CEO was forced to reveal that he received a bonus of $1623,753 for 1929, while six vice presidents received $1.4 million in aggregate.[37] The revelations—coming at the beginning of the Great Depression—sparked a variety of shareholder lawsuits demanding that the executives return up to $36.5 million in bonuses received since 1911. The same year, shareholders sued American Tobacco for details on its stock subscription plan, resulting in revelations that the company's CEO netted $1.2 million from an incentive plan that allowed him to purchase company stock at deeply discounted prices.[38] Wells (2010, p. 712) concludes that "the Bethlehem Steel and American Tobacco revelations, combined no doubt with a Depression-generated disgust with corporate management, fueled public perceptions that executive compensation was both excessive and the product of self-dealing."

---

[37] "$1,623,753 Grace's Bonus For 1929: Bethlehem President Testifies At Merger Trial To Receiving This Amount," *Wall Street Journal* (1903), "Bonus Figures Given At Trial: Six Vice Presidents Of Bethlehem Received $1,432,033 In 1929", *Wall Street Journal* (1930).

[38] In particular, American Tobacco's George Hill was allowed to purchase 13,440 shares of company stock at its $25 "par value" at a time when shares were trading for about $120. See "G. W. Hill Got Bonus of $1,200,000 Stock", *New York Times*(1931).

## 3.3  Depression-Era Outrage and Disclosure Requirements (1930s)

We have become accustomed to the idea that shareholders—and the public in general—have a right to know the details of the compensation paid to top executives in publicly traded corporations. However, the initial push for pay disclosure was not driven by shareholders but rather by "New Deal" politicians outraged by perceived excesses in executive compensation.

In 1933, Franklin D. Roosevelt became president, ending three terms and twelve years of Republican government and ushering in the New Deal in a country recovering from the Great Depression. In the April prior to the 1932 election—in the face of proposed bailout loans from the government's Reconstruction Finance Corporation (RFC)—the Interstate Commerce Commission demanded that all railroads disclose the names of executives making more than $10,000 per year.[39] The disclosed pay levels outraged the new Administration, and in May 1933 the RFC required railroad companies receiving government assistance to reduce executive pay by up to 60%.[40] Ultimately, the US Senate authorized the Federal Coordinator of Transportation to impose an informal (but uniformly complied-with) cap of $60,000 per year for all railroad presidents.

The mandated pay disclosures for railroad executives sparked the interest of other US regulators. By mid-1933 the Federal Reserve began investigating executive pay in its member banks, the RFC conducted a similar investigation for non-member banks, and the Power Commission investigated pay practices at public utilities. In October 1933, the Federal Trade Commission (FTC) requested disclosure of salaries and bonuses paid by all corporations with capital and assets over $1 million (approximately 2000 corporations).[41] Business leaders questioned whether the FTC had the legal authority to compel such disclosures, but were reminded that, "Congress in its present temper would readily authorize" whatever the FTC wanted.[42] Executives were particularly incensed that the FTC would demand such closely guarded information without any explanation of how the information would be used and without any confidentiality guarantees.

---

[39] "Railroad Salary Report: I.C.C. Asks Class 1 Roads About Jobs Paying More Than $10,000 a Year," *Wall Street Journal* (1932).

[40] The required reductions ranged from 15% (for executives earning less than $15,000) to 60% (for executives earning more than $100,000. See "RFC Fixed Pay Limits: Cuts Required to Obtain Loans," *Los Angeles Times* (1933), "Cut High Salaries or Get No Loans, is RFC Warning", *New York Times* (1933).

[41] See Robbins, "Inquiry into High Salaries Pressed by the Government", *New York Times* (1933) and "President Studies High Salary Curb: Tax Power is Urged as Means of Controlling Stipends in Big Industries", *New York Times* (1933). In addition to investigating corporate executive pay, President Roosevelt personally called attention to lavish rewards in Hollywood, resulting in a provision added to the moving-picture code that imposed heavy fines on companies paying unreasonable salaries.

[42] "Federal Bureau Asks Salaries of Big Companies' Executives", *Chicago Daily Tribune* (1933).

Following the Securities Act of 1934, the responsibility for enforcing pay disclosures for top executives in publicly traded corporations was consolidated into the newly created Securities and Exchange Commission (SEC). In December 1934, the SEC issued permanent rules demanding that companies disclose the name and all compensation (including salaries, bonuses, stock, and stock options) received by the three highest-paid executives. The securities of companies not complying with the new regulations by June 1935 would be removed from exchanges. Several companies, including US Steel, pleaded unsuccessfully for the SEC to keep the data confidential, arguing that publication "would be conducive to disturbing the morale of the organization and detrimental to the best interests of the registrant and its stockholders".[43]

Under the Securities Act, details on executive pay are disclosed in company proxy statements issued in connection with the company's annual shareholders meeting. Ultimately, these disclosures have provided the fodder for all subsequent pay controversies. Proxy statements for companies with December fiscal closings are typically issued in late March or early April, triggering a deluge of pay-related articles in the popular and business press each Spring. *Forbes* and *Business Week* began offering extensive lists of the highest-paid executives in 1970. *Fortune* and the *Wall Street Journal* quickly followed suit, and by now most major newspapers conduct their own CEO pay surveys for companies based in their local metropolitan areas.

While the SEC has no direct power to regulate the level and structure of CEO pay, the agency *does* determine what elements of pay are disclosed and how they are disclosed. The SEC has routinely expanded disclosure requirements from year to year, with major overhauls in 1978, 1992, 2006, and 2011. The first proxy statements issued after the formation of the SEC were typically about three-to-five pages long, with less than one page devoted to executive compensation. By 2007, the average proxy statement exceeded 70 pages, nearly all focused on compensation.[44]

Under the theory that sunlight is the best disinfectant, the SECs disclosure rules have long been a favorite method used by the SEC and Congress in attempts to curb perceived abuses and excesses in executive compensation. Indeed, most additions to disclosure requirements over time—including perquisite disclosure in the 1970s, enhanced option grant disclosures in the 1992, and actuarial pension values in 2006—reflect policy responses to relatively isolated abuses. However, there is little evidence that enhanced disclosure leads to reductions in objectionable practices: for example, perquisites increased as executives learned what was common at other firms, and options exploded following the 1993 rules.

---

[43] "US Steel Guards Data on Salaries: Sends details confidentially to SEC head with request that they be kept secret", *New York Times* (1935).

[44] The average length of 2007 proxy statements for the 100 largest firms (ranked by revenues) was 62.8 pages (ignoring appendices). In 2006—before the 2006 disclosure rules—the average length was 45 pages.

The demand for disclosure reflects both legitimate shareholder concerns and public curiosity. While disclosure can conceptually facilitate better monitoring of outside directors by shareholders, the public curiosity aspect of disclosure imposes large costs on organizations. The recurring populist revolts against CEO pay, for example, could not have been waged without public pay disclosure. Public disclosure effectively ensures that executive contracts in publicly held corporations are not a private matter between employers and employees but are rather influenced by the media, labor unions, and by political forces operating inside and outside companies. These "uninvited guests" to the bargaining table have no real stake in the companies being managed and no real interest in seeing companies managed well so they serve all the claimants on the firm including consumers, debt and equity holders, employees and communities. However, as will become evident throughout this section, these third parties have affected both the level and structure of executive pay through tax policies, accounting rules, direct legislation, and other rules and regulations stretching back nearly a century. These important but often ignored costs of disclosure must be weighed against the benefits (better monitoring of directors) in determining the optimal amount of pay disclosure for top managers.

## 3.4 The Rise (and Fall) of Restricted Stock Options (1950–1969)

In the 1920s, the US income tax was new, the use of stock options was new, and no one had figured out whether options would be taxed: (1) as compensation when options are exercised (and hence taxed as ordinary income for the individual, and representing a deductible business expense for the company); or (2) as capital gains when the stock purchased upon exercise was ultimately sold (and hence taxed at a lower capital gains rate for the individual, with the company forgoing deductibility). It took nearly twenty years for this issue to be resolved. The major case study at the time involved a May 1928 option grant to the CEO of a chain of movie theaters. After a large six-month run-up in the stock price following the grant, the CEO exercised his options in October 1928 and subsequently sold the shares in 1929 and 1930, paying capital gains taxes (12.5%) on the proceeds. The Bureau of Internal Revenue (the predecessor of the Internal Revenue Service (IRS)) held that he owed ordinary income taxes on the spread at exercise (25% in 1928). The taxpayer appealed the decision, and nearly nine years later the Circuit Court of Appeals agreed with the taxpayer, concluding that a taxable gain is realized only when the shares are sold and not when the option is exercised.[45] However, the Bureau appealed this decision, and in a related case nine years later, ruled in favor of the Bureau, concluding in 1946 that the gain upon exercise is compensation, thereby taxable as ordinary income.[46]

[45] Rossheim v. Commissioner, 92 F. 2d 247 (1937).
[46] Commissioner v. Smith, 324 U.S. 177 (1945).

By 1950, the tax issue surrounding stock options was a big deal: the highest marginal tax rate on ordinary income and corporate profits had swelled to 91% and 50.75% (from 25% and 12% in 1928, respectively), compared to a capital gains rate of 25% (from 12.5% in 1928). Moreover, while the Supreme Court required taxes to be paid immediately upon exercise, the 1934 Securities Act required executives to hold shares acquired through option exercises for at least six months before they could sell.[47] For example, suppose an executive acquired one share of stock at an exercise price of $10 when the market price is $25. To finance the exercise and pay the taxes, the executive would need to pay $23.65 (i.e. the exercise price plus 91% of the exercise-date spread), but could not raise the amount by selling shares.

As part of the Revenue Act of 1950, a business-friendly Congress unhappy with the recent Supreme Court decision created a new type of stock options called "restricted stock options" that would be taxable not upon exercise but only when the shares were ultimately sold (and then taxed as capital gains). Restricted stock options solved the tax-timing problem, since taxes were not owed until the stock was sold (at least six months following the exercise date). Given the tax rates at the time, restricted stock options also became a relatively efficient way to convey after-tax compensation to executives. For example, at a 91% tax rate on ordinary income and 50.75% corporate tax rate, it cost shareholders $5.47 in after-tax profit to give the executive $1 in after-tax income.[48] In contrast (and for simplicity ignoring the timing issues), when the pay is taxed as capital gains rather than ordinary income, it cost shareholders only $1.33 to convey $1 in after-tax income to the executive (even though shareholders forfeit the deduction).

The passage of the 1950 Act launched a predictable wave of new option plans. In 1950 approximately 4% of the companies listed on the New York Stock Exchange (NYSE) had option plans for their top executives; by June 1951 the number had tripled to 12%.[49] In their study of the fifty largest firms in 1940 and 1960, Frydman and Saks (2010) estimate that the fraction of executives holding stock options increased from less than 10% in 1950 to over 60% by 1960. Grant sizes also grew: the grant–date value of options for those executives receiving options increased from about 10% of total compensation in the early 1950s to over 20% of total compensation by the early 1960s.

---

[47] To deter insider trading, Section 16b of the 1934 Securities Act requires that any profit realized by an officer or director in the purchase or sale of an equity security within a six-month period be returned to the company.

[48] At a 91% tax rate, the CEO must receive $11.11 before tax to realize $1 after tax. But, at a 50.75% corporate tax rate, paying $11.11 in deductible compensation costs reduces after-tax profits by only $5.47.

[49] Mullaney, "Parley Here Indicates the Continued Spread in Industry of Stock Purchase Option Plans," *New York Times* (1951). The percentages are based on average of 840 NYSE-listed firms in 1950 and 876 in June 1951.

Figure 13 shows the average level and structure of compensation for CEOs in 50 large manufacturing firms, based on data from Lewellen (1968).[50] The stock option data—compiled long before the availability of option-pricing methodologies such as Black and Scholes (1973)—are based on appreciations in the annual spread between the market and exercise prices of outstanding options. Since Lewellen measures options at their appreciated values, the trend in Figure 13 reflects, in part, general stock-market movements over this time period. After adjusting for inflation, salaries and bonuses fell from over $2.2 million in 1940 to about $1.5 million (in 2011-constant dollars) from 1947 to 1963. Total compensation, including deferred compensation and stock options, peaked at $2.9 million in 1956. Negligible before 1951, options grew to over 30% of compensation by 1956, falling to about a fifth of total compensation by 1963.

Since restricted stock options were taxed at a much lower rate than salaries, the trend in Figure 13 understates the growing importance of options on an after-tax basis. In particular, Lewellen estimates that options accounted for nearly half of total after-tax compensation in 1956, falling to a third of total after-tax compensation by 1963.

By the summer of 1951, there was a growing backlash against the perceived escalation in restricted stock option plans. In August 1951, the Salary Stabilization Board conducted a series of hearings on whether stock options should be considered compensation under the Defense Production Act and therefore subject to regulation by the Stabilization Board.[51] In November 1951, the Stabilization Board ruled that restricted stock options could be granted without the Board's approval as long as the option met certain conditions (including an exercise price of at least 95% of the grant-date stock price; restricted options with an exercise price as low as 85% of the stock price could be issued, but would be considered increases in salary subject to regulation).[52] The Board's ruling was followed by a second wave of option plans, and by June 1952 nearly 17% of the NYSE firms had adopted plans.[53] In July 1952 the Salary Stabilization Board was disbanded.

---

[50] Lewellen (1968) reports both the pre-tax and after-tax values for salaries and bonuses, but only the after-tax values for stock options and deferred compensation. The pre-tax values for stock options after 1950 are determined by dividing the after-tax value by .85 (Lewellen uses a 15% effective tax rate for options). The pre-tax value for deferred compensation (and for options prior to 1950) are estimated by dividing the after-tax value by (1-t★), where t★ is one-half of the implied average tax rate for salaries and bonuses. For example, if Lewellen reports pre-tax and after-tax salaries and bonuses of $240,000 and $80,000, respectively, suggesting an average tax rate of 60%, we would calculate pre-tax deferred compensation using a tax rate of 30%.

[51] "Salary Board's Panel to Study Stock Option in Top Executive Pay," *Wall Street Journal* (1951), "Options Defended at Salary Hearing: Restricted Stock Plans Called Neither Inflationary Nor Compensatory by 8 Men," *New York Times* (1951), "Options on Stocks Scored at Hearing: Majority of Witnesses Call it Inflationary and Unfair to Small Stockholders," *New York Times* (1951), "Salary Board Urged to Ban Stock Option Plans Until End of Emergency," *Wall Street Journal* (1951), "Stock Options: Industry Says Salary Board Should Keep Its Hands Off Employee Plans," *Wall Street Journal* (1951).

[52] "Rules are Issued on Stock Options," *New York Times* (1951).

[53] "One in 6 Companies Gives Stock Options," *New York Times* (1952).

**Figure 13** Trends in before-tax CEO Compensation in 50 Large Manufacturing Companies, 1940–1963. *Note:* The figure is based on the Lewellen (1968) study of 50 large manufacturing firms, adjusted for inflation using the Consumer Price Index. The value of stock options is based annual calculations of the spread between the market and exercise prices. The before-tax value of deferred pay and stock options are estimated from Lewellens after-tax calculations.

Many of the options granted in the early 1950s fell underwater in the 1953 post–Korean War recession. As part of the Revenue Act of 1954, Congress modified the restrictions on restricted stock options by officially sanctioning variable–price options, in which the exercise price of a previously granted option could be lowered if it turned out that the market price of the optioned stock declined subsequent to the granting of the option. In addition, where the 1950 Act put no limits on the expiration terms of options, the 1954 Act limited exercise terms to 10 years (which continues to be the most common term for options granted through current times). While the popularity of stock options decreased briefly during the bear market in 1957,[54] the use of stock options continued to trend upward: by 1961, 68% of the NYSE firms had option plans.[55]

During the 1960 recession, as new option grants were falling out of favor given the declining stock market, companies began exploiting the provision of the 1954 Act allowing repricing of options by either resetting exercise prices or by canceling existing options and replacing them with new options with lower exercise prices. This practice

[54] "Ailing Options: Stock Market Decline Dulls Allure of Plans For Company Officials," *Wall Street Journal* (1957).

[55] The 1961 survey is described in Stanton, "Cash Comeback: Stock Options Begin to Lose Favor in Wake of Tax Law Revision," *Wall Street Journal* (1964).

became highly controversial in the early years of the Kennedy Administration, leading to a series of Congressional hearings aimed at repealing the favorable tax treatment for restricted stock options.[56] In 1961, the President demanded that Congress remove the favorable tax treatment for options, instead taxing options as ordinary income upon exercise (most of which would be subject to the 91% top marginal tax rate). The issue was debated in Congress for the next two years, and the controversy intensified in late 1963 and early 1964 when it was revealed that executives at Chrysler had realized $4.2 million in gains from exercising stock options in 1963, and had sold nearly 200,000 shares acquired through earlier exercises.[57] Ultimately, as part of the Revenue Act of 1964, Congress stopped short of removing the favorable tax status of restricted stock options, but took several steps that substantially reduced their attractiveness. In particular, under the new law:

- Executives were required to hold stock acquired through option exercises for three years (rather than six months) in order to be taxed at the lower capital gains rate.
- Exercise prices could be no less than 100% (rather than 85%) of the grant-date market prices.
- The maximum option term was reduced from ten years to five years.
- The option price could not be reduced during the term of the option, nor could an option be exercised while there is an outstanding option issued to the executive at an earlier time. (This provision was designed to halt the practice of repricing options or canceling out-of-the-money options and replacing them with options with lower exercise prices).

To distinguish options meeting these new requirements from restricted options granted under the Revenue Act of 1950 provisions, the 1964 Act referred to new grants as "qualified stock options" rather than restricted stock options.

Finally (but perhaps most importantly), the 1964 law reduced the top marginal tax rate on ordinary income from 91% to 70%, which significantly reduced the attractiveness of restricted options over cash compensation. Figure 14 provides a historical comparison of the tax advantages of restricted or qualified stock options relative to cash compensation or non-qualified stock options (in which the gains upon exercise are taxed as ordinary income for the recipient, and deductible as a compensation expense to the company). As a result of the 1964 tax law, the after-tax cost to investors of conveying

[56] "Options on the Wane: Fewer Firms Plan Sale of Stock to Executives at Fixed Exercise Prices," *Wall Street Journal* (1960), "Congress and Taxes: Specialists Mull Ways to Close "Loopholes" in Present Tax Laws," *Wall Street Journal* (1959), "House Group Hears Conflicting Views on Stock Option Taxes," *Wall Street Journal* (1959).

[57] "Chrysler Chairman Defends Option Plan, Offers to Discuss It With Federal Officials," *Wall Street Journal* (1963), "Chrysler Officers Got Profit of $4.2 Million On Option Stock in '63," *Wall Street Journal* (1964), "Chrysler Officers' Sale of Option Stock Could Stir Tax Bill Debate," *Wall Street Journal* (1963), "House Unit Seen Favoring Curbs on Stock Options," *Wall Street Journal* (1963), "Senate Unit Votes to Tighten Rules on Stock Options," (1964).

**Figure 14** The Tax Acts of 1964 and 1969 reduced the tax advantages of restricted or qualified stock options Note: The figure shows the after-tax cost to investors of conveying an incremental $1 in after-tax income under two tax regimes: (1) ordinary compensation (taxable to the recipient at the top marginal rate for earned income ($t_I$), and deductible by the firm at the top marginal rate for corporate income ($t_C$)), and (2) capital gains (taxable to the recipient at the capital gains rate ($t_G$), but not deductible by the firm). The cost for ordinary income is computed as $(1-t_C)/(1-t_I)$, while the cost for capital gains is $1/(1-t_G)$.

an after–tax dollar to the CEO in cash compensation fell from $5.56 to $1.73, while the cost of conveying an after–tax dollar in restricted or qualified stock options (taxed as capital gains) remained at $1.33.

The popularity of qualified stock options fell as a result of the 1964 tax law[58] and collapsed following the Tax Reform Act of 1969. In addition, the 1969 Act defined gains from exercising restricted or qualified options as a tax preference item subject to a new Alternative Minimum Tax (AMT) on high wage earners.[59] The 1969 Act gradually reduced the top marginal tax rate on earned income from 77% in 1969 to 50% by 1972, reduced the corporate tax rate from 52.8% to 48%, and raised the top capital gains tax

---

[58] See Stanton, "Cash Comeback: Stock Options Begin to Lose Favor in Wake of Tax Law Revision," *Wall Street Journal* (1964). Stock options briefly resurged in 1966, following at 25% increase in the Dow Jones average from 1964 to early 1966 (Elia, "Opting for Options: Stock Plans Continue in Widespread Favor Despite Tax Changes," *Wall Street Journal* (1967)).

[59] In particular, if the option gains exceed 50% of the executive's total income (including option gains), the amount of the option gain over 50% would be treated as fully taxable ordinary income. The AMT was passed following revelations that 155 high–income households took deductions that reduced their federal tax liabilities to zero.

rate from 25% in 1969 to 36.5%. Once the new rates were fully implemented (and ignoring AMT issues), it cost investors approximately $1.04 in after-tax profit to convey an incremental $1 in after-tax income to the CEO through cash compensation or non-qualified stock options, and $1.57 to convey $1 in qualified stock options. Thus, for executives and companies in the highest tax brackets, qualified stock options became tax disadvantageous compared to non-qualified stock options, and (as illustrated in Figure 14) have remained so throughout the early 2000s. Indeed, Hite and Long (1982) provide evidence that the 1969 Act explains the dramatic shift from qualified stock options to non-qualified stock options that took place during the early 1970s. Restricted or qualified stock options—which had been the dominant form of long-term incentives for two decades—virtually disappeared.

## 3.5  Wage-and-Price Controls and Economic Stagnation (1970–1982)

### 3.5.1  America, Land of the Freeze

In August 1971, in an ultimately (and predictably) unsuccessful attempt to control inflation, President Nixon imposed a 90-day freeze on commodity prices and wages (including executive pay). In December 1971—in what was called Phase Two of the Nixon wage-and-price controls—the Pay Board established by Congress imposed a limit of 5.5% for increases in executive pay (the limit being binding for company-defined groups of executives, but not necessarily for individual executives).[60] The Nixon wage-and-price controls were not the first time that levels of executive compensation were explicitly limited by legislation, but were the first time such controls were imposed in a peacetime economy. In particular, the World War II-era Stabilization Act of 1942 froze wages and salaries (for executives as well as other labor groups) at their September 15, 1942 level. The Stabilization Act expired in 1946, but was replaced during the Korean War by the Salary Stabilization Boards established in May 1951 as part of the Defense Production Act of 1951. Similar to the Nixon controls, the Korean War Salary Board set a 6% limit on pay increases for each company's executives taken as a group; the limits were lifted when the Board was quietly disbanded in July 1952.[61]

In a debate (and outcome) eerily similar to what would happen two decades later during the Clinton Administration, concerns that the Nixon wage controls would significantly reduce executives incentives led to a series of compromises (or loopholes, depending on one's perspective).[62] In particular, while bonuses were generally limited to the amount paid in any one of the last three years plus 5.5%, the limit did not apply

---

[60] Hunt, "Board Agrees on Tightening of Standards on Executive Pay, Increases Topping 5.5%," *Wall Street Journal* (1971).

[61] "Old Wage Board Exits: New Unit to Take Over with Reduced Powers," *Wall Street Journal* (1952).

[62] Jensen, "Bonuses Rise Through Loopholes," *New York Times* (1972). For the complete text of the executive compensation provisions, see "Board's Text on Executive Compensation," *Wall Street Journal* (1971).

to existing sales incentives, commission and production-incentive programs. Moreover, companies could petition to adopt new incentive plans as long as they were directly related to increased productivity. As a result, scores of companies introduced performance-based bonus plans tied to accounting data or revenues, or converted their existing plans into plans exempt from the limits.

Non-qualified stock options were allowed under the Nixon controls only if the plan was shareholder-approved, if the aggregate number of options granted did not increase from the prior three years, and if the exercise price was at least 100% of the grant-date market price. Non-qualified options were treated as wages and salaries under the Nixon controls, and were valued at 25% of the fair-market value of the shares underlying the option.[63] This valuation approach represents an interesting (albeit short-lived) historical footnote, since it was imposed a year before Black and Scholes (1973) and decades before companies began routinely placing a value on options when making compensation decisions.

The median continuing CEO in the *Forbes 800* received a 4.5% increase in cash compensation in 1971 (below the Nixon limit), 6.0% in 1972, and 8.1% in 1973 (both above the Nixon limit).[64] Since the government-mandated limits on pay raises applied only to executives taken as a group and not individual executives, companies routinely raised CEO pay by reducing pay (or offering smaller raises) to lower-level executives.[65] In August 1973, to stop companies from raising CEO pay above the 5.5% limit, the Nixon Administration imposed the 5.5% limit on the more-narrowly defined group of executives identified in company proxy statements. The wage-and-price controls expired in May 1974, in spite of Administration efforts to retain limits on executive compensation.[66]

CEO pay rose significantly after the wage controls were lifted in May 1974. The median continuing CEO in the *Forbes 800* received an 11.1% increase in nominal cash compensation in 1974, double the average limit under the Nixon controls. From 1973 through 1979, the median cash compensation for CEOs in the *Forbes 800* increased by 12.2% each year (doubling from $162,000 to $324,000), significantly exceeding the average annual inflation rate of 8.5%.

Figure 15 shows the median level and structure of compensation for CEOs in 73 large manufacturing firms from 1964 to 1982, based on data from Murphy (1985) and

[63] Valuation is based on testimony by Richard McNamar, director of the Pay Boards office of economic policy. See Calame, "Executives' Pay Faces Going-Over By Wage Board," *Wall Street Journal* (1972).

[64] The calculations are based on annual compensation surveys published in *Forbes* covering the largest 500 companies ranked by revenues, assets, market capitalization, and employees (about 800 companies are listed in one or more of these *Forbes* rankings annually).

[65] "Government Moves to Hold Executives to 5.5% Pay Boosts," *Wall Street Journal* (1973).

[66] "Business Groups Oppose Nixon Control Plan, Intensify Their Efforts to Abolish Restraints," *Wall Street Journal* (1974), "Nixon Halts Push to Retain Some of Phase 4 Controls," *Wall Street Journal* (1974).

**Figure 15** Trends in before-tax CEO Compensation in 73 Large Manufacturing Companies, 1964–1982. *Note:* The figure is based on data from the Murphy (1985) study of executive pay in 73 large manufacturing firms, adjusted to 2011 dollars using the Consumer Price Index. Stock options are valued on grant-date using the Black and Scholes (1973) formula.

inflation–adjusted to 2010–constant dollars. To my knowledge, this was the first compre-hensive study of executive pay that measured stock options as the grant–date value using the Black and Scholes (1973) approach. In nominal terms (that is, before adjusting for inflation), median CEO pay in the 73 firms in Figure 15 nearly tripled from $148,900 in 1964 to $569,550 in 1982. However, after adjusting for inflation (which averaged over 6.5% annually over this period), median real CEO pay increased only by about 23% over this 18–year period, or about 1.2% per year. Stock options accounted for 2% of total pay for the average CEO in 1964; the use of options had grown to 12% of pay by 1981. Both the level of pay and the use of stock options fell during the 1981-1982 recession.

### 3.5.2 The Controversy over Perquisites

While cash compensation escalated (at least in nominal terms) during the 1970s, the use of stock options was relatively stagnant. Part of the declining popularity of options reflected the change in tax policies in 1964 and 1969 that made qualified stock options less attractive, coupled with their outright prohibition in 1976 (see below). More importantly, though, was the prolonged stagnation in the stock market, driven in part by the oil–price shocks of 1973 and 1977. In particular, the nominal value of the bell–wether Dow Jones average was basically flat from the beginning of 1965 through the early 1980s (falling from 903 in January 1965 to below 800 by mid–January 1982, and

only surpassing 1050 on one day over these seventeen years). While executives continued to receive periodic option grants during this time (once every three years was typical), many of the grants replaced options that expired worthless or options that were cancelled and reissued with a lower exercise price.

The void in compensation created by worthless stock options was quickly filled by a plethora of new plans designed to provide more predictable payouts, including: book-value plans (where executives receive dividends plus the appreciation in book values); long-term performance plans (with payouts based on long-term earnings growth targets); and guaranteed bonuses (with payouts guaranteed independent of performance).[67] In addition, since the Nixon wage-and-price controls restricted salaries but not company-provided benefits, companies began relying to a greater extent on shareholder-subsidized perquisites or perks such as low-interest loans, yachts, limousines, corporate jets, club memberships, hunting lodges and corporate retreats at exotic locations.

By the mid-1970s, perceived abuses attracted the ire of shareholder activists, the SEC and the IRS.[68] In December 1975, the IRS circulated a draft of proposed regulations specifying which fringe benefits could be excluded from an executive's taxable income. A long-held general rule excluded from taxable income benefits arising from the ordinary course of business that do not cost the employer anything extra (such as family members accompanying an executive on the corporate jet). The proposed rule imposed tax liabilities for these and other fringe benefits if the benefits were available only to the most highly compensated executives.

The attack on perquisites escalated in 1977 as President Carter famously rallied against companies taking deductions for the three-martini lunch, yachts and hunting lodges maintained to entertain business associates, first-class air travel, fees paid to social and athletic clubs and money spent on sports and theater tickets.[69] Congress resisted implementing most of Carter's reforms as part the Revenue Act of 1978 (in large part because it would potentially affect their *own* consumption of perquisites) but agreed to eliminate deductions for entertainment facilities.[70]

[67] Ricklefs, "Sweetening the Pot: Stock Options Allure Fades, So Firms Seek Different Incentives," *Wall Street Journal* (1975), Hyatt, "No Strings: Firms Lure Executives By Promising Bonuses Not Linked To Profits," *Wall Street Journal* (1975), Ricklefs, "Firms Offer Packages of Long-Term Incentives as Stock Options Go Sour for Some Executives," *Wall Street Journal* (1977).

[68] See, for example, Bender, "The Executive's Tax-Free Perks: The IRS Looks Harder at the Array of Extras," *New York Times* (1975a), Bender, "Fringe Benefits at the Top: Shareholder Ire Focuses on Loan Systems," *New York Times* (1975b), Blumenthal, "Misuse of Corporate Jets by Executives is Drawing More Fire," *New York Times* (1977), Schellhardt, "Perilous Perks: Those Business Payoffs Didn't All Go Abroad; Bosses Got Some, Too; IRS and SEC Investigating Loans and Lush Amenities Provided for Executitves; An Eye on Hunting Lodges," *Wall Street Journal* (1977).

[69] Rankin, "Incentives for Business Spending Proposed in Corporate Package," *New York Times* (1978), "Excerpts From Carter Message to Congress on Proposals to Change Tax System," *New York Times* (1978).

[70] Zimmerman, "Washington Word: Don't Do as We Do But Do as We Say: For Bureaucrats, Lawmakers, Hard Times Aren't Here; Limousines and Free Trips," *Wall Street Journal* (1975).

In August 1977, the SEC issued Interpretive Release #5856 stating that the value of perquisites be included as compensation in proxy statements.[71] In justifying the new disclosures, SEC enforcement chief Stanley Sporkin argued that the "excesses just got to the point where it became a scandal".[72] The disclosures in the 1978 proxy statements fueled the fire by focusing even more attention on perquisites.[73] The information on perquisites was expanded significantly in 1979 proxy statements, when the SEC implemented its first major revision in proxy disclosures since the 1930s. Also in 1979, the IRS issued significant new auditing guidelines aimed at detecting and taxing executive perquisites. McGahran (1988) argues that the new SEC disclosures made it easier for the IRS to detect (and tax) fringe benefits, and presents some evidence that fringe benefits decreased, while cash compensation increased, as a result of the SEC and IRS actions.

The ongoing attack on perquisites was reflected in the contemporaneous early academic literature on agency theory. For example, the "agency problem" introduced by Jensen and Meckling (1976) focused on managerial consumption of non-pecuniary benefits such as "the physical appointments of the office" and "the attractiveness of the secretarial staff". Similarly, Alchian and Demsetz (1972) conclude that companies allow personal consumption of corporate (or university) property (such as "privileges, perquisites, or fringe benefits") because the cost of detecting and punishing such "turpitudinal peccadilloes" is larger than the benefits from prohibiting the activity.

### 3.5.3  There's No Accounting for Options

The restricted and qualified stock options created by the 1950 and 1964 Revenue Acts were not formally considered compensation and therefore companies did not record an expense for such options for either tax or accounting purposes. The switch to non-qualified options in the 1970s—which were considered compensation for tax purposes—raised a new question: how should options be accounted for in company income statements? One possibility was to follow the tax code by recognizing an accounting expense at the time an option is exercised. But, in spite of its simplicity, this method is inconsistent with the basic tenet of accounting that expenses should be matched to the time period when the services associated with those expenses were rendered. Rather, the tenet suggested that options should be expensed over their term based on the grant-date value of the option. At the time, however (and for a long time to come) there was no accepted way of placing a value on an employee stock option.

---

[71] "Personal-Use Perks For Top Executives Are Termed Income: SEC Says Valuable Privileges Will Have to be Reported As Compensation by Firms," *Wall Street Journal* (1977).

[72] Jensen, "Executives' Use of Perquisites Draws Scrutiny," *New York Times* (1978).

[73] Examples include: Joseph, " US Industries Faces Queries on its Perks At Annual Meeting," *Wall Street Journal* (1978), Metz, "Close Look Expected At Executive Perks in Proxy Material: SEC Stress on Disclosure Is Linked to Coming Tales of Holder-Assisted Goodies," *Wall Streety Journal* (1978), Penn, "Ford Motor Covered Upkeep for Elegant Co-Op of Chairman: Questions Arise on Personal vs. Business Use of Suite in Posh New York Hotel," *Wall Street Journal* (1978).

In October 1972, the Accounting Principles Board (APB)—the predecessor to the current Financial Accounting Standards Board (FASB)—issued APB Opinion No. 25, "Accounting for Stock Issued to Employees". Under APB Opinion No. 25, the compensation expense associated with stock options was defined as the (positive) difference between the stock price and the exercise price as of the first date when both the number of options granted and the exercise price become known or fixed. The expense for this spread between the price and exercise price—called the intrinsic value—was amortized over the period in which the employee is prohibited from exercising the option.[74] Under this rule, there was no charge for options granted with an exercise price equal to (or exceeding) the grant-date market price, because the spread is zero on the grant-date.

The accounting treatment of options cemented the dominance of the traditional stock option (an option granted with a five- or ten-year term with an exercise price equal to the grant-date market price) and discouraged companies from offering more novel option plans. For example, APB Opinion 25 imposes a higher accounting charge for options with an exercise price indexed to the stock-price performance of the market or industry, because the exercise price is not immediately fixed. Similarly, it imposes a higher accounting charge for options that only become exercisable if certain performance triggers are achieved, because the number of options is not immediately fixed. Finally, it imposes an accounting charge for options that are issued in the money but not for options issued at the money—a feature that became especially significant three decades later in the scandals involving backdating.

### 3.5.4 The Rise (and Fall) of Stock Appreciation Rights

Under Section 16(b) of the Securities Act of 1934, executives must return any profits realized from buying and selling (or selling and buying) shares of their company's stock within any period of less than six months. This constraint was not problematic for executives exercising restricted or qualified stock options, since the provisions of the 1951 and 1964 Revenue Acts already required executives to hold shares for six months (for restricted options) or three years (for qualified options) before trading. However, the six-month holding period was particularly troublesome for non-qualified options, since executives were required to pay ordinary income tax based on the date the option was exercised and not when the underlying shares were sold.[75] Given the depressed stock market in the 1970s, the taxes due upon exercise were often greater than value of the shares when they became tradable.

---

[74] This period is often called the vesting period but this terminology is misleading since vesting implies that the executive is free to sell the option or keep it if he leaves the firm, as opposed to being able only to exercise the option.

[75] The executive could defer the taxes during the six-month holding period, but would still owe taxes on the gain on the exercise date even if stock prices fell over the subsequent six months.

In December 1976, the SEC formally exempted stock appreciation rights (SARs) from the Section 16(b) short-swing profit prohibition.[76] Executives holding a SAR are entitled to receive the appreciation on one share of stock. Like stock options, SARs had a pre-determined term but executives were generally free to exercise their SARs at any time prior to the end of this term (after some minimum time had elapsed). Prior to the December 1976 ruling, there was considerable debate about whether SARs would be subject to the short-swing rule and therefore the proceeds from the exercise of the rights would have to be returned to the company. After the SEC ruling, SARs provided a way for executives to reap the benefits of exercising non-qualified options without being subject to the six-month holding requirement.[77] As a result of the ruling, many companies replaced their option grants with SAR grants, or issued tandem SARs and options, which allowed the executive to decide which to exercise. For the next fifteen years, SARs became a ubiquitous component of long-term compensation for most executives.

Jumping ahead a bit, in May 1991 the SEC declared that the six-month holding period begins when options are granted, and not when executives acquire shares upon exercise. Therefore, as long as the executive has held the option for at least six months, he is allowed to immediately sell the shares acquired when options are exercised. This new ruling eliminated the primary advantage of SARs over non-qualified options and, as a result, SARs largely disappeared from existence. In addition, the SEC rule effectively encouraged the practice—commonplace today—of selling shares immediately upon exercise.[78]

The rise and ultimate fall of SARs is a tribute to the cleverness of companies in finding ways around rules that disadvantage executives and companies (in this case, the six-month holding requirement).[79] Moreover, the experience shows how seemingly innocuous government interventions (in this case, the 1976 and 1991 SEC rulings) can have a major impact on the composition of executive compensation.

### 3.5.5  Qualified Stock Options Resurrected, But No One Cares

The Revenue Acts of 1964 and 1969 significantly reduced the attractiveness of restricted/qualified stock options, but did not prohibit new grants. As part of the Revenue Act of 1976, Congress allowed executives to retain and exercise grants made

---

[76] "SEC Exempts Rights To Stock Appreciation From 'Insider' Curbs," *Wall Street Journal* (1976).

[77] There was one major disadvantage of SARs over non-qualified options: companies granting SARs were required to record an accounting charge for the evolving value of the SARs, while there was typically no accounting charge for options.

[78] Peers, "Executives Take Advantage of New Rules on Selling Shares Bought With Options," *Wall Street Journal* (1991).

[79] A related innovation in the late 1980s was the "Stock Depreciation Right," which provided cash payments to executives exercising options if stock prices fell during the six month holding period. (See Crystal, "The Wacky, Wacky World of CEO Pay," (1988)).

prior to May 20, 1976, but banned all future grants of qualified stock options. Since existing grants had a maximum five-year term, the last grant of qualified options was set to expire on May 19, 1981.

As 1981 approached, Congress resurrected a new form of qualified options (now called Incentive Stock Options or ISOs) as a last-minute addition to the Economic Recovery Tax Act of 1981.[80] ISOs carried many of the restrictions common for qualified stock options (holding periods after exercise, minimum exercise prices, etc.), and in addition were limited to $100,000 per executive per year (calculated as the stock price multiplied by the number of options on the date of grant). While ISOs have continued to be popular in the 2000s for middle-level managers (where the $100,000 limitation is not binding) and for companies without taxable profits (where loss of deductibility for ISOs is not costly), virtually all options granted to CEOs and other top executives since 1972 have been non-qualified stock options.

### 3.5.6  Bigger is Better (Paid)

Almost half of the cross-sectional variation in cash compensation in the United States between 1970 and 1982 was explained by company size (usually measured by firm revenues), and the highest-paid executives routinely were at the helm of the largest conglomerates and largest steel, automotive, and oil companies. Year-to-year changes in cash compensation were also largely driven by increases in company size. And non-monetary aspects of compensation—including power, prestige, board memberships, and community standing—were also positively linked to increases in firm size. The strong relation between CEO pay and company size gave CEOs substantial incentives to increase company size, while the decline of equity-based incentive plans gave them little incentive to increase company share prices. It is noteworthy that the implicit incentives to increase company revenue help explain the unproductive diversification, expansion, and investment programs in the 1970s, which in turn further depressed company share prices.

Although CEO pay and bottom-line corporate profitability remained relatively stagnant from 1970 to 1982, productivity did not. Spurred in part by the oil-price shocks of 1973 and 1977, this period brought significant technological advances that improved productivity, declines in regulation, and increases in global trade significant enough to constitute what Jensen (1993) calls the "Modern Industrial Revolution". By the early 1980s, most sectors in the US economy were saddled with increasing excess capacity, implying that the sectors had more capital and labor than were required to maintain current levels of production. The root causes of the excess capacity differed across industries. In the oil sector, for example, the five-fold increase in the inflation-adjusted price of crude oil led firms to launch massive capacity-increasing exploration

---

[80] Bettner, "Incentive Stock Options Get Mixed Reviews, Despite the Tax Break They Offer Executives," *Wall Street Journal* (1981).

and development projects in anticipation of continued price increases; the sector was stuck with the capacity when demand dropped and prices tumbled to pre-shock levels. Technological change dramatically increased capacity for computing firms, while increased competition from non-unionized entrants created excess capacity in a variety of industries ranging from steel to groceries.

By definition, investment in an industry with excess capacity is a negative net-present-value project, since the industry already has more capital and labor than can be productively employed. Indeed, firms with excess capacity can either increase output with the same workforce, or maintain current output with a smaller workforce. However, the 1970s conglomerates and other large companies typically chose to neither increase output (given low market demands) nor decrease their workforce (since pecuniary and non-pecuniary rewards for CEOs were both tied to company size). Moreover, by the end of the 1970s, most of these companies were generating huge amounts of cash, far in excess of that required to fund available positive net present value projects. CEOs, loathe to distribute excess cash back to shareholders, responded by wasting huge amounts of free cash flow through unwise diversification and investment programs.[81]

## 3.6  The Emerging Market for Corporate Control (1983–1992)

### 3.6.1  Golden Parachutes and Section 280(G)

The executive compensation practices of the 1970s provided few incentives for executives to pursue value-increasing reductions in excess capacity and disgorgements of excess cash. Equity-based compensation (mostly in the form of stock options) accounted for only a small fraction of CEO pay (Figure 15) and the options that existed often were underwater or expired worthless. Annual bonuses—the dominant form of compensation-based incentives—were focused on beating annual budget targets rather than creating long-run value. Performance-based terminations were almost non-existent and—since the vast majority of CEO openings were filled by incumbents rather than outside hires—the managerial labor market was similarly ineffective in disciplining poor performance.

Boards of directors, typically dominated by corporate insiders (in influence if not in numbers), had little reason to reduce corporate waste as long as the companies were delivering positive nominal profits. However, pressures to improve performance and disgorge cash were ultimately introduced by the capital markets, including "hostile takeover" artists such as Carl Icahn, Irwin Jacobs, Carl Lindner, David Murdock, Victor Posner, Charles Bluhdorn, and T. Boone Pickens. At the time, these takeover artists were known pejoratively as "corporate raiders", though history has shown they were a

---

[81] Jensen (1986a) defines free cash flow as cash flow in excess of that which can be reinvested at returns equal to or better than the cost of capital.

positive force in creating substantial amounts of value for shareholders of target firms while reallocating resources to higher-valued uses.[82] Sometimes this wealth was created by the post-merger activities of the raiders (such as firing incompetent incumbent managers and selling non-productive assets). At other times the wealth was created by responses to the takeover threat (such as spending cash to repurchase shares or to purchase competitors, causing resources to leave the sector and allowing shareholders to find more productive uses for their cash).

The takeover market was complemented by the emergence of leveraged buyouts (LBOs): going-private transactions financed by debt using the target firm's future cash flows as collateral. Debt created value by providing commitments that the firm would pay its cash flows to debtholders, reducing the amounts available for executives to waste (Jensen, 1986a). Debt also taught executives that capital is costly (since the interest cost of debt capital was more obvious than the implicit, though larger and largely unrecognized, cost of equity capital), leading to reductions in inventories and working capital. The emergence of LBOs and leveraged recapitalizations (in which the firm leverages the capital structure while staying public) created substantial amounts of shareholder value in firms with stable cash flows and no productive alternative uses for their cash, characteristics of many of the mature and declining sectors in the early 1980s.

While employment in companies targeted by hostile takeovers or LBOs was modestly reduced (which was productive given the presumptive excess capacity), the individuals most vulnerable to job losses were incumbent executives opposed to the changes in control. Innovations designed to thwart takeovers included greenmail payments (repurchase of the raiders' stock at above market prices), standstill agreements (bribes so that the raider does not purchase additional stock), staggered boards (where directors serve overlapping terms, making it difficult for a proxy fight to gain a majority), supermajority rules (requiring more than 50% votes to approve a merger), and poison pills (where shareholders get special rights when there is a takeover bid). But, perhaps the most notorious innovation was the "golden parachute" which provided direct payments to executives following a successful change in control. In most cases, the payment required both the change-of-control and the loss of a job (hence, called "double-triggered" since two things had to happen); in other cases (single-triggered) the change-of-control itself was sufficient to trigger the payment, regardless of job loss.[83]

---

[82] See Holderness and Sheehan (1985) for an analysis of how the first six on this list improved operating results and shareholder values, and Fischel (1995) for an analysis of how T. Boone Pickens facilitated the restructuring of the oil sector.

[83] In regulations associated with the TARP bailouts in 2008–2009, Congress redefined golden parachutes to refer to any severance payment in connection with an executive departure, regardless of whether the departure was related to a change of control. In contrast, the golden parachute label prior to the TARP bailouts required a change of control, but did not require departure. For example, accelerated vesting of restricted stock or accelerated exercisability of stock options upon a change of control was considered part of the parachute payment, even if the executive retained his or her job.

Whether change-of-control agreements facilitate or thwart takeovers remains a matter of debate and rests in the details. On one hand, as emphasized by Jensen (1986b), such agreements facilitate takeovers by providing bribes to existing managers to acquiesce to the change in control. On the other hand, such agreements can significantly increase the cost of takeovers for prospective acquirers, especially if the agreements cover dozens or hundreds of executives who have no plausible influence over the takeover decision. In any case, the existence of the apparent bribes paid to top executives (but not to shareholders in general) attracted the ire of a Congress already skeptical of hostile takeovers and their benefits.

Change-in-control arrangements became controversial following a $4.1 million payment to William Agee, the CEO of Bendix. In 1982, Bendix launched a hostile takeover bid for Martin Marietta, which in turn made a hostile takeover bid for Bendix. Bendix ultimately found a "white knight" and was acquired by Allied Corp., but only after paying CEO Agee the golden parachute. The payment sparked outrage in Washington, but Congress could not ban golden parachute payments outright, because such a ban would pre-empt state corporation laws. Congress does, however, control the tax laws, which allow corporations to deduct compensation from income only if the payments represent reasonable compensation for services rendered. By defining particular types or dollar amounts of compensation as unreasonable, Congress can directly determine whether compensation is deductible for corporate tax purposes.

Congress attempted to discourage golden parachutes by adding Sections 280(G) and 4999 to the tax code as part of the Deficit Reduction Act of 1984. Section 280(G) of the Code provides that, if change-in-control payments exceed three times the individuals base amount, then *all* payments in excess of the base amount are non-deductible to the employer. Also, Section 4999 imposes a 20% excise tax on the recipient of a parachute payment on the amount of payment above the base amount. The base amount is typically calculated as the individuals average total taxable compensation (i.e. W-2 compensation, which include gains from exercising stock options) paid by the company over the prior five years.

Because of the complexity of what appears to be a simple rule, modest increases in parachute payments can trigger substantial tax payments by both the company and executive. For example, suppose an executive with five-year average taxable compensation of $1 million receives a golden parachute payment of $2.9 million, which is less than three times the $1 million base amount.[84] In this case, the entire $2.9 million parachute payment would be deductible by the company, and would be taxable as ordinary income to the executive. In contrast, suppose that the golden parachute payment

---

[84] The golden parachute payment includes not only cash payments but also the value of accelerated vesting of stock and options, as long as the payment is contingent on a change of control or ownership of the company.

was $3.1 million, which is more than three times the $1 million base amount. Under Section 280(G), the company would not be able to deduct $2.1 million (of the $3.1 million parachute payment) as a compensation expense, and (under Section 4999) the executive would owe $420,000 in excise taxes (i.e. 20% of $2.1 million) in addition to ordinary income taxes on the full $3.1 million parachute payment.

The new Section 280(G) impacted executive compensation in several ways. First, the new law led to a proliferation in change-in-control agreements, which had previously been fairly rare. The Deficit Reduction Act was signed into law on July 18, 1984. By 1987, 41% of the largest 1000 corporations had golden parachute agreements for their top executives, and the prevalence of golden parachutes increased to 57% in 1995 and to 70% by 1999.[85] In addition, the standard golden parachute payment quickly became the government prescribed amount of three times base compensation. By 1991, 47.5% of CEO golden parachute arrangements specified a multiple of three times base pay, and by 1999 71% specified three times base pay. Thus, the rule designed to limit the generosity of parachute payments has led to both a proliferation and a standardization of Golden Parachute payments in most large corporations. Apparently compensation committees and executives took the regulation as effectively endorsing such change-in-control agreements as well as the payments of three times average compensation (which quickly became the standard).

Second, Section 280(G) (and the corresponding Section 4999) gave rise to the "tax gross up", in which the company offset the tax burden of the 20% excise tax by paying an additional amount for the tax (and the tax on the additional amount).[86] The percentage of agreements that included gross-up provisions increased from 38% in 1991 to over 82% by 1999.[87] This gross-up concept was subsequently applied to a variety of executive benefits with imputed income taxable to the executive, such as company cars, club memberships, and personal use of corporate aircraft.

Third, Section 280(G) also provided incentives for companies to shorten vesting periods in stock option plans, and incentives for executives to exercise stock options even earlier than they would normally be exercised. Consider two otherwise identical executives with golden parachutes paying three times base compensation and holding identical options. Suppose that one of the executives exercises a year prior to the change in control, while the other holds until the change in control. Since base compensation under Section 280(G) includes gains from exercising options, the first executive can receive a higher parachute payment before triggering the excise tax, thus increasing the

---

[85] Alpern and McGowan (2001), p. 6.

[86] For example, continuing with the example above, suppose the CEO owed $420,000 in excise taxes (i.e. 20% of the $2.1 million excess benefit). If the CEO had a gross-up clause (and assuming a marginal tax on ordinary income of 50% on top of the 20% excise tax), he would receive a gross-up payment of $1.4 million and a total change-in-control payment of $4.5 million, leaving him with after-tax income of $1.55 million (which is what he would have received without an excise tax).

[87] Alpern and McGowan (2001, p. 7–8).

benefits from early exercise. Moreover, unexercisable stock options routinely become vested (or exercisable) upon a change in control, and the value of these options is defined by the IRS as part of the parachute payment subject to the excise taxes. Therefore, companies and executives can reduce change-in-control related tax liabilities by shortening the time until options become exercisable, and by exercising early and therefore reducing the incentive effects of those plans.

Similarly, unvested restricted stock routinely becomes vested upon a change in control, and a portion of the value of these shares upon vesting is defined by the IRS as part of the parachute payment subject to the excise taxes. Thus, companies can also reduce change-in-control related tax liabilities by shortening the vesting period for restricted stock.

Finally, but perhaps most importantly, the 1984 tax laws regarding Golden Parachutes appear to have triggered the proliferation of Employment Agreements for CEOs and other top-level executives in most large firms since the mid-1980s. In particular, Section 280(G) applies only to severance payments contractually tied to changes of control, while individual CEO employment agreements typically provide for severance payments for *all* forms of terminations without cause, including (but not limited to) terminations following control changes. Therefore, companies can circumvent the Section 280(G) three-times-base-compensation limitations (at a potentially huge cost to shareholders) by making payments available to all terminated executives, and not only those terminated following a change in control. Indeed, Graef Crystal (when he was still a leading compensation consultant) predicted the unintended consequences of the enactment of these tax provisions in his 1984 opinion piece in the *Wall Street Journal:*

> But will Congress's new reforms really curb those who want to offer excessive compensation? Not necessarily. Congress has, as usual, made an opening move in a corporate chess game and neglected to consider its opponents countermoves. Instead of having a contract that covers only a change of control, some companies may now implement all-embracing employment contracts that guarantee a person employment (or what he would have earned had he continued to be employed), for say, five years, and under all circumstances. You won't see one word in that contract about payments in the event of a change of control, and the net effect will be to give the executive more than he would have had had Congress not given free rein to its passions.[88]

In summary, although Section 280(G) was meant to reduce the generosity of parachute payments, the government action appears to have increased the prevalence of: (i) change-in control plans; (ii) tax gross-ups; (iii) early exercise of stock options; (iv) short vesting periods for restricted stock and stock options; and (v) employment agreements. Each of these outcomes both reduces the incentive effects of incentive compensation for CEOs and other executives and increases the costs of these plans to their firms.

---

[88] See Crystal, "Manager's Journal: Congress Thinks It Knows Best About Executive Compensation," *Wall Street Journal* (1984).

### 3.6.2 The Shareholder Awakening

The emerging market for corporate control had pronounced effects on the US stock market. After nearly two decades of stagnation, the Dow Jones Industrial Average rallied from below 800 to over 2700 between mid-1982 and mid-1987 (i.e. appreciating nearly 30% per year for five years). While the largest beneficiaries were shareholders in firms that became takeover targets, the rally was broad based and lifted share prices across a wide range of firms and industries. However, executives vigorously (and often successfully) fought takeovers in the 1980s by adopting anti-takeover provisions and by lobbying for political protection (Holmstrom and Kaplan, 2001). Therefore, in spite of the gains to shareholders (or perhaps because of the redistribution in wealth resulting from these gains), hundreds of bills were introduced in Congress to curb takeovers and highly leveraged transactions (Fischel, 1995).

Court decisions and legislation in the late 1980s (coupled with the October 1987 stock market crash) brought the hostile takeover market in the US to a virtual halt. The high-yield debt market was crippled by the indictment and subsequent guilty pleas of Michael Milken and Drexel Burnham Lambert and by restrictions on high-yield debt holdings imposed on savings institutions, commercial banks, and insurance firms, and by major punitive changes in the US bankruptcy law that made it uneconomic to reorganize troubled firms outside of bankruptcy.

But the lessons of the wealth creations learned from the takeover wave resonated with shareholders. In 1985, Robert Monks founded Institutional Shareholder Services to provide proxy-voting advice to institutional shareholders. In 1986, corporate raider T. Boone Pickens founded the United Shareholders Association focused on improving governance and compensation. Academics increasingly argued that traditional management incentives that focused on company size, stability, and accounting profitability destroyed rather than created value, and recommended that executive pay be tied more closely to company value through increases in stock options and other forms of equity-based incentives. These pressures began having an impact: non-equity-based CEO pay continued to grow in real terms after the mid-1980s, but became a smaller part of the total compensation package. For the first time since the 1950s, stock options re-emerged as the dominant form of incentives compensation.

Figure 16 shows the median level and average structure of CEO compensation from 1980-1992, based on Hall and Liebman (1998). Total grant-date compensation is defined as the sum of salaries, bonuses, and the grant-date value of stock options using the Black and Scholes (1973) formula. The annual sample size varies between 365 and 432 firms, and is representative of the population of the large US firms. The percentage composition is defined by dividing the average salary and bonus (or options) by the average total compensation for each year.[89] As shown in the figure, inflation-adjusted

---

[89] The percentage compositions in Figure 16 are not strictly comparable to those in Figure 4 or Figure 15, and overstate the percentage of compensation in options relative to the methodology used elsewhere in this study.

**Figure 16** Median Grant-date Compensation for CEOs in Hall and Liebman (1998), 1980–1992. *Note:* The figure is based on data from the Hall and Liebman (1998) study of executive pay in approximately 500 large US firms (the annual sample size varies between 365 and 432 firms). Total compensation, adjusted to 2011 dollars using the Consumer Price Index, includes salaries, bonuses, and stock options valued on the grant-date using the Black and Scholes (1973) formula. Pay composition percentages for each year based on the annual sample averages for the two components.

median pay levels doubled from 1980 to 1992 from $946,000 to $1,800,000. The increase in pay primarily reflects the increase in stock option grants, which accounted for nearly half of total aggregate CEO pay by 1992.

Although the takeover and LBO market had been largely shut down by political forces, investors and executives began recognizing that value is created through reducing excess capacity or by reversing ill-advised diversification programs. As emphasized by Holmstrom and Kaplan (2001), stock options allowed executives to share in the value created by internal restructurings: "Shareholder value became an ally and not an enemy".

### 3.6.3 Controversial Pay Leads to Sweeping New Disclosure Rules

Between October 13–19, 1987, the Dow Jones Average dropped nearly 800 points (from 2508 to 1738), losing 30% of its value in a week. Executive stock options, which had only recently become an important part of pay, were suddenly underwater. Companies responded by repricing existing options or by significantly increasing the size of their post-crash option grants (Saly, 1994).

The October 1987 crash turned out to be short-lived: by August 1989 the Dow Jones reached an all-time high of 2735, hitting 3000 by July 1990. Stock options issued

both before and after the crash were well in the money and becoming exercisable. Large manufacturing firms—still sorting out the excess capacity issues of the 1970s—were downsizing and laying off workers, to the delight of shareholders but attracting the ire of Congress, labor unions, and the media. The combination of valuable options, robust stock markets, and the 1990-1991 recession provided the perfect recipe for a populist attack on executive pay.

The CEO pay debate achieved international prominence during the 1990-1991 recession. The controversy heightened with the November 1991 introduction of Graef Crystal's (1991) exposé on CEO pay, *In Search of Excess*, and exploded following President George H. W. Bush's pilgrimage to Japan in January 1992 (accompanied by an entourage of highly paid US executives). What was meant to be a plea for Japanese trade concessions dissolved into accusations that US competitiveness was hindered by its excessive executive compensation practices as attention focused on the huge pay disparities between top executives in the two countries.[90]

In response to growing outrage, legislation was introduced in the House of Representatives disallowing deductions for compensation exceeding 25 times the lowest-paid worker, and the Corporate Pay Responsibility Act was introduced in the Senate to give shareholders more rights to propose compensation-related policies. The SEC pre-empted the pending Senate bill in February 1992 by requiring companies to include non-binding shareholder resolutions about CEO pay in company proxy statements,[91] and announced sweeping new rules in October 1992 affecting the disclosure of top executive compensation in the annual proxy statement. Among other changes, the SEC's new 1992 disclosure rules required companies to produce (a) a Summary Compensation Table summarizing the major components of compensation received by the CEO and other highly paid executives over the past three years; (b) tables describing option grants, option holdings, and option exercises in greater detail; (c) a chart showing the company's stock-price performance relative to the performance of the market and their peer group over the prior five fiscal years; and (d) a report by the compensation committee describing the company's compensation philosophy. Overall, the new rules dramatically increased the information available about stock option grants and holdings, and the performance graph cemented the idea that the objective of the firm was to create shareholder value.

## 3.7  The Stock Option Explosion (1992–2001)

As shown in Figure 17 (and Figures 3 and 4), the median pay for CEOs in S&P 500 firms more than tripled between 1992 and 2001, driven by an explosion in the use of stock options. CEO incentive compensation in the early 1990s was split about evenly

---

[90] "SEC to Push for Data on Pay of Executives," *Wall Street Journal* (1992).
[91] "Shareholder Groups Cheer SEC's Moves on Disclosure of Executive Compensation," *Wall Street Journal* (1992).

**Figure 17**  Median Grant-date Compensation for CEOs in S&P 500 Firms, 1992–2001. *Note:* Median pay levels (in 2011-constant dollars) based on ExecuComp data for S&P 500 CEOs. Total compensation (indicated by bar height) defined as the sum of salaries, non-equity incentives (including bonuses), benefits, stock options (valued on the date of grant using company fair-market valuations, when available, and otherwise using ExecuComps modified Black–Scholes approach), stock grants, and other compensation. Other compensation excludes pension-related expenses. Pay composition percentages defined as the average composition across executives.

between options and accounting-based bonuses. By 1996, options had become the largest single component of CEO compensation in S&P 500 firms, and the use of options was even greater in smaller firms (and especially high-tech start-ups). By 2000, stock options accounted for more than half of total compensation for a typical S&P 500 CEO.

The escalation of stock-option grants cannot be explained by a single factor. Instead, I believe that there are six main factors that fueled the explosion in stock options:

- *Shareholder pressure for equity-based pay;*
- *SEC holding-period rules;*
- *Clinton's $1 million deductibility cap;*
- *Accounting rules for options;*
- *SEC option disclosure rules;*
- *NYSE listing requirements.*

In this section, I will discuss each of these factors in rough chronological order (referring to prior discussions when appropriate), and indicate how they contributed to the option explosion.

### 3.7.1  Shareholder Pressure for Equity-Based Pay

As discussed in Section 3.6.2, the decline in takeover activity in the late 1980s corresponded to the rise in shareholder activism. This new breed of activists—including many of the largest state pension funds—demanded increased links between CEO pay and shareholder returns. The activists were joined by academics such as Jensen and Murphy (1990a), who famously (or infamously) argued "It's not *how much* you pay, but *how* that matters." Jensen and Murphy (1990b) showed that CEOs of large companies were paid like bureaucrats, in the sense that they were primarily paid for increasing the size of their organizations, received small rewards for superior performance, even smaller penalties for failures, and that the bonus components of the pay packages showed very little variability, less even then the variability of the pay of rank-and-file employees. They concluded that compensation committees and boards should focus primarily on the incentives provided by the pay package rather than the level of pay, and were joined by shareholder activists such as the United Shareholders Association in advocating more stock ownership and more extensive use of stock options.

Companies responded by taking Jensen and Murphy's mantra a bit too literally: adding increasingly generous grants of stock options on top of already competitive pay packages, without any reduction in other forms of pay and showing little concern about the resulting inflation in pay levels.

### 3.7.2  SEC Holding–Period Rules

When an executive exercises a non-qualified stock option, the executive pays the exercise price and owes income tax on the gain. As discussed in Section 3.5.4, SEC rules in effect May 1991 required executives to hold shares acquired from exercising stock options for at least six months. The executive could defer the taxes during the six-month holding period (leading many executives to exercise after June 30, pushing the tax liability to the following year), but would still owe taxes on the gain on the exercise date even if stock prices fell over the subsequent six months. This rule implies that executives cannot finance the exercise by selling shares acquired in the exercise, and executives exercising stock options therefore faced both significant short-run cash–flow problems (from paying the exercise price) and increased risk.

Before May 1991, the SEC defined the exercise of an option as a "stock purchase" reportable by corporate insiders on Form 4 within 10 days following the month of the transaction. On May 1, 1991, in response to demands for more transparency of option grants, the SEC defined the *acquisition* rather than the exercise of the option as the reportable stock purchase. As a consequence of this change, the six-month holding period required by the Securities Act's "short-swing profit" rule now begins when options are granted, and not when executives acquire shares upon exercise. Therefore, as long as the options are exercised more than six months after they are granted, the

executive is free to sell shares immediately upon exercise. This ruling significantly increased the value of the option from the standpoint of the recipient.

### 3.7.3  The Clinton $1 Million Deductibility Cap

The controversy over CEO pay became a major political issue during the 1992 US presidential campaign.[92] Bill Clinton promised to end the practice of allowing companies to take unlimited tax deductions for excessive executive pay; Dan Quayle warned that corporate boards should curtail some of the exorbitant salaries paid to corporate executives that were unrelated to productivity; Bob Kerry called it unacceptable for corporate executives to make millions of dollars while their companies were posting losses; Paul Tsongas argued that excessive pay was hurting America's ability to compete in the international market; and Pat Buchanan argued "you can't have executives running around making $4 million while their workers are being laid off."

After the 1992 election, president-elect Clinton reiterated his promise to define compensation above $1 million as unreasonable, thereby disallowing deductions for all compensation above this level for all employees. Concerns about the loss of deductibility contributed to an unprecedented rush to exercise options before the end of the 1992 calendar year, as companies urged their employees to exercise their options while the company could still deduct the gain from the exercise as a compensation expense.[93] In anticipation of the loss of deductibility, large investment banks accelerated their 1992 bonuses so that they would be paid in 1992 rather in 1993. In addition, several publicly traded Wall Street firms, including Merrill Lynch, Morgan Stanley, and Bear Stearns, announced that they consider returning to a private partnership structure if Clinton's plan were implemented.[94]

By February 1993, President Clinton backtracked on the idea of making *all* compensation above $1 million unreasonable and therefore non-deductible, suggesting that exemptions would be granted if the company could meet (not yet developed) federal standards proving that the executive improved the firm's productivity.[95] In April, details of the considerably softened plan began to emerge.[96] As proposed by the Treasury Department and eventually approved by Congress as part of the Omnibus Budget Reconciliation Act of 1993, Section 162(m) of the tax code applies only to public firms and not to privately held firms, and applies only to compensation paid to the CEO and

---

[92] "Politics and Policy—Campaign '92: From Quayle to Clinton, Politicians are Pouncing on the Hot Issue of Top Executive's Hefty Salaries," *Wall Street Journal* (1992).

[93] Chronicle Staff and Wire Reports, "Big Earners cashing in now: fearful of Clinton's tax plans, they rush to exercise their options," *San Francisco Chronicle* (1992).

[94] Siconolfi, "Wall Street is upset by Clinton's support on ending tax break for 'excessive' pay," *Wall Street Journal* (1992).

[95] Freudenheim, "Experts see tax curbs on executives' pay as more political than fiscal," *New York Times* (1993), Ostroff, "Clinton's Economic Plan Hits Taxes, Payrolls and Perks," (1993).

[96] Greenhouse, "Deduction proposal is softened," *New York Times* (1993).

the four highest-paid executive officers as disclosed in annual proxy statements (compensation for all others in the firm is fully deductible, even if in excess of the million-dollar limit). More importantly, Section 162(m) does not apply to compensation considered performance-based for the CEO and the four highest-paid people in the firm.

Performance-based compensation, as defined under Section 162(m), includes commissions and pay based on the attainment of one or more performance goals, but only if (1) the goals are determined by an independent compensation committee consisting of two or more outside directors, and (2) the terms of the contract (including goals) are disclosed to shareholders and approved by shareholders before payment. Stock options generally qualify as performance based, but only if the exercise price is no lower than the market price on the date of grant. Base salaries, restricted stock vesting only with time, and options issued with an exercise price below the grant-date market price do not qualify as performance based.

Under the IRS definition, a bonus based on formula-driven objective performance measures is considered performance based (so long as the bonus plan has been approved by shareholders), while a discretionary bonus based on ex post subjective assessments is not considered performance based (because there are not predetermined performance goals). However, the tax law has been interpreted as allowing negative but not positive discretionary payments: the board can use its discretion to pay less but not more than the amount indicated by a shareholder-approved objective plan.

In enacting Section 162(m), Congress used (or abused) the tax system to target a small group of individuals (the five highest-paid executives in publicly traded firms) and to punish shareholders of companies who pay high salaries. Indeed, the stated objective of the proposal that evolved into Section 162(m) was not to increase tax revenues but rather to reduce the level of CEO pay. For example, the House Ways and Means Committee described the congressional intention behind the legislation:

> Recently, the amount of compensation received by corporate executives has been the subject of scrutiny and criticism. The committee believes that excessive compensation will be reduced if the deduction for compensation (other than performance-based compensation) paid to the top executives of publicly held corporations is limited to $1 million per year.[97]

Ironically, although the objective of the new IRS Section 162(m) was to reduce excessive CEO pay levels by limiting deductibility, the ultimate result (similar to what happened in response to the golden parachute restrictions) was a significant *increase* in CEO pay. First, since compensation associated with stock options is generally considered performance-based and therefore deductible (as long as the exercise price is at or above the grant-date market price), Section 162(m) encouraged companies to grant more stock options. Second, while there is some evidence that companies paying base salaries in excess of $1 million lowered salaries to $1 million following the enactment of

[97] 1993 US Code Congressional and Administrative News 877, as cited in Perry and Zenner (2001).

Section 162(m) (Perry and Zenner, 2001), many others raised salaries that were below $1 million to exactly $1 million (Rose and Wolfram, 2002). Finally, companies subject to Section 162(m) typically modified bonus plans by replacing sensible discretionary plans with overly generous formulas (Murphy and Oyer, 2004).

It is difficult to argue with the principle that companies should only be able to deduct compensation expenses for services rendered. However, the $1 million reasonableness standard is inherently arbitrary and has not been indexed for either inflation (+60% from 1993 to 2011) or changes in the market for executive talent: compensation plans that seemed excessive in 1993 are considered modest by current standards. More importantly, Section 162(m) disallows deductions for many value-increasing plan designs. For example, Section 162(m) disallows deductions for restricted stock or for options issued in the money, even when such grants are accompanied by an explicit reduction in base salaries. In addition, Section 162(m) disallows deductions for discretionary bonuses based on a board's subjective assessment of value creation. I suspect that many compensation committees have welcomed the tax-related justification for not incorporating subjective assessments in executive reward systems. After all, no one likes receiving unfavorable performance evaluations, and few directors enjoy giving them. It is therefore not surprising that directors are often unwilling to devote the time, personal effort and courage to provide accurate, frank and effective performance appraisals of CEOs and other top executives. But, by failing to make the appraisals, directors are breaching one of their most important duties to the firm.

Moreover, Section 162(m) has distorted the information companies give to shareholders. In particular, in order to circumvent restrictions on discretionary bonuses, companies have created a formal shareholder-approved plan that qualifies under the IRS Section 162(m) while actually awarding bonuses under a different shadow plan that pays less than the maximum allowed under the shareholder-approved plan. These shadow plans often have little or nothing to do with the performance criteria specified in the shareholder-approved plans. As a consequence, the bonus plans and the performance metrics described in company proxy statements are not necessarily reflective of the actual formulas and performance measures used to determine bonuses.

Finally, it is worth noting that Section 162(m) is highly discriminatory, applying only to the compensation received by the top five executive officers, and applying only to publicly traded companies and not to private firms or partnerships. Ultimately, arbitrary and discriminatory tax rules such as Section 162(m) have increased the cost imposed on publicly traded corporations and have made going-private conversions more attractive.

### 3.7.4  There's (Still) No Accounting for Options

The 1972 APB Opinion 25—which defined the accounting treatment for stock options as the spread between the market and exercise price on the grant-date—pre-dated Black and Scholes (1973), which offered the first formula for computing the value of a traded

call option. Academic research in option valuation exploded over the next decade, and financial economists and accountants became increasingly intrigued with using these new methodologies to value, and account for, options issued to corporate executives and employees.

In 1984, the Financial Accounting Standards Board (FASB) floated the idea that companies account for employee stock options using the so-called minimum value approach.[98] By June 1986, the FASB idea had evolved into a proposal with the important change that the accounting charge would be based on the fair market value (e.g. the Black–Scholes value) and not a minimum value. The proposal was vehemently opposed by all of the Big Eight accounting firms, the American Electronics Association (including more than 2800 corporate members), the Financial Executives Institute, the Pharmaceutical Manufacturers Association, and the National Venture Capital Association.[99] Many of the criticisms focused on the complexity of the Black–Scholes formula, as exemplified by the following quote from Joseph E. Connor, chairman of Price Waterhouse:

> *Corporate America rightfully is skeptical of any standard that depends upon complex pricing models that provide partial and debatable answers. Yet after two years of fruitless efforts, the FASB persists in trying to turn this ivory-tower notion into a usable standard. The compensation element is a mirage, tempting to the imagination but impossible to touch. The board should turn its attention to more productive areas.[100]*

Ultimately, and without fanfare, FASB tabled its 1986 proposal before submitting an exposure draft.

In late 1991, Senator Carl Levin (D–Michigan) attempted to bypass FASB by introducing the Corporate Pay Responsibility Act requiring companies to take a charge to their earnings to reflect the cost of option compensation packages; as noted in Section 3.6.3, the bill also directed the SEC to require more disclosure for stock option arrangements in company proxy statements. Although Levin's bill was ultimately shelved, it provided pressure for renewed FASB focus on option expensing.

---

[98] The minimum value approach is identical to the value of a forward contract to purchase a share of stock at some date in the future at a pre-determined price (that is, an option without the option to refrain from buying when the price falls below the exercise price). For example, the minimum value of an option on a non-dividend-paying stock is calculated as the current stock price minus the grant-date present value of the exercise price. Thus, the value of a ten-year option granted with an exercise price of $30 when the grant-date market price was $25 would be $V = \$25 - \$30/(1+r)^{10}$, where r is the risk-free rate.

[99] See, for example, Rudnitsky and Green, "Options Are Free, Aren't They?", *Forbes*, (August 26, 1985), Gupta and Berton, "Start-up Firms Fear Change in Accounting," *Wall Street Journal* (1986), Gupta and Berton, "Start-up Firms Fear Change in Accounting," *Wall Street Journal* (1986), Fisher, "Option Proposal Criticized," *New York Times* (1986), Eckhouse, "Tech Firms' Study: Accounting Rule Attacked," *San Francisco Chronicle* (1987).

[100] Connor, "There's No Accounting for Realism at the FASB," *Wall Street Journal* (1987).

In April 1992, FASB voted 7-0 to endorse an accounting charge for options, and issued a formal proposal in 1993. The proposal created a storm of criticism among business executives, high-tech companies, accountants, compensation consultants, the Secretary of the Treasury, and shareholder groups.[101] In March 1994, more than 4000 employees from 150 Silicon Valley firms rallied against the accounting change, calling on the Clinton Administration to block the proposal because it would restrict job creation and economic growth. Even President Clinton, usually a critic of high executive pay, waded into the debate by expressing that it would be unfortunate if FASB's proposal inadvertently undermined the competitiveness of some of America's most-promising high-tech companies.[102] In the aftermath of the overwhelmingly negative response, FASB announced it was delaying the proposed accounting change by at least a year, and in December 1994 it dropped the proposal.

In 1995, FASB issued a compromise rule, FAS 123, which *recommended* but did *not require* that companies expense the fair-market value of options granted (using Black–Scholes or a similar valuation methodology). However, while FASB allowed firms to continue reporting under APB Opinion 25, it imposed the additional requirement that the value of the option grant would be disclosed in a footnote to the financial statements. Predictably, only a handful of companies adopted FASB's recommended approach. As I will discuss below in Section 3.8.4, it wasn't until the accounting scandals in the early 2000s that a large number of firms voluntarily began to expense their option grants.

The accounting treatment of options promulgated the mistaken belief that options could be granted without any cost to the company. This view was wrong, of course, because the opportunity or economic cost of granting an option is the amount the company could have received if it sold the option in an open market instead of giving it to employees. Nonetheless, the idea that options were free (or at least cheap) was erroneously accepted in too many boardrooms. Options were particularly attractive in cash-poor start-ups (such as in the emerging new economy firms in the early 1990s), which could compensate employees through options without spending any cash. Indeed, providing compensation through options allowed the companies to generate cash, since when options were exercised, the company received the exercise price and could also deduct the difference between the market price and exercise price from its corporate taxes. The difference between the accounting and tax treatment gave option-granting companies the best of both worlds: no accounting expense on the company's books, but a large deduction for tax purposes. When coupled with the May 1991 rule

---

[101] See, for example, Berton, "Business chiefs try to derail proposal on stock options," *Wall Street Journal* (1992), Harlan and Berton, "Accounting Firms, Investors Criticize Proposal on Executives' Stock Options," *Wall Street Journal* (1992), "Bentsen Opposes FASB On Reporting Stock Options," *Wall Street Journal* (1993), Berton, "Accounting Rule-Making Board's Proposal Draws Fire," *Wall Street Journal* (1994), Harlan, "High Anxiety: Accounting Proposal Stirs Unusual Uproar In Executive Suites," *Wall Street Journal* (1994).
[102] "Clinton Enters Debate Over How Companies Reckon Stock Options," *Wall Street Journal* (1993).

eliminating holding requirements after exercise, stock options had important perceived advantages over all other forms of compensation.

As both an illustration of how accounting affects compensation decisions, and as an interesting episode in its own right, consider how a change in accounting rules affected option repricing. On December 4, 1998, FASB announced that repriced options issued on or after December 15, 1998 would be treated under "variable accounting", meaning that the company would take an accounting charge each year for the repriced option based on the actual appreciation in the value of the option. FASB issued its final rule in March 2000 as FASB Interpretation No. 44, or FIN 44, indicating that FASB did not consider this a new rule but rather a re-interpretation of an old rule. In particular, FASB reasoned that the "fixed accounting" under APB Opinion 25 (in which the option expense was equal to the spread between the market and exercise price on the first date when both the number of options granted and the exercise price become known or fixed) did not apply to companies that have a policy of revising the exercise price.

Companies with underwater options rushed to reprice those options in the 12-day window between December 4–15, 1998.[103] Indeed, Carter and Lynch (2003) document a dramatic increase in repricing activities during the short window, followed by dramatic declines; Murphy (2003) shows that repricings virtually disappeared after the accounting charge. Many companies with declining stock prices circumvented the accounting charge on repriced options by canceling existing options and re-issuing an equal number of options after waiting six months or more. But this replacement is not neutral. It imposes substantial risk on risk-averse employees since the exercise price is not known for six months and can conceivably be *above* the original exercise price. In addition, canceling and reissuing stock options in this way provides perverse incentives to keep the stock-price down for six months so that the new options will have a low exercise price. All of this scrambling to avoid an accounting charge!

### 3.7.5  SEC Option Disclosure Rules

The most widely debated issue surrounding the SEC's 1992 disclosure rules was how stock options would be valued in both the Summary Compensation Table and in the Option Grant table. The SEC wanted a total dollar cost of option grants so that the components in the Summary Compensation Table could be added together to yield a value for total compensation, and lobbied for calculating option cost using a Black and Scholes (1973) or related approach. The SEC's proposal was vehemently opposed by high option-granting firms (especially from the Silicon Valley and Boston's 128 corridor) and (more surprisingly) by compensation consulting and accounting firms. Ultimately, a compromise was struck: the Summary Compensation Table would include

---

[103] Johnston, "Fast Deadline On Options Repricing: As of Next Tuesday, It's Ruled an Expense," *New York Times* (1998).

the number, but not the cost, of options granted, thus defeating the SEC's objective of reporting a single number for total compensation. In addition, companies would have a choice in the Option Grant Table to report either the Black–Scholes grant-date cost or the *potential* cost of options granted (under the assumption that stock prices grow at 5% or 10% annually during the term of the option).[104]

From the perspective of many boards and top executives who perceive options to be nearly costless, or indeed deny that options have value when granted, the only way they can quantify the options they award is by the number of options granted. The focus on the quantity rather than the cost of options is further solidified by the SEC's 1992 disclosure rule and also by institutions that monitor option plans. For example, under the current listing requirements of the New York Stock Exchange and the National Association of Security Dealers (NASD), companies must obtain shareholder approval for the total number of options available to be granted, but not for the cost of options to be granted. Advisory firms (such as Institutional Shareholder Services) often base their shareholder voting recommendations on the option "overhang" (that is, the number of options granted plus options remaining to be granted as a percent of total shares outstanding), and not on the opportunity cost of the proposed plan. Therefore, boards and top executives often implicitly admit that the *number* of options granted imposes a cost on the company, while at the same time denying that these options have any real dollar cost to the company.

The focus on the quantity rather than the cost of options granted helps explain a puzzling result in the executive pay literature (e.g. Hall and Murphy, 2003): the near-perfect correlation between the S&P 500 Index and average grant-date CEO pay. Figure 18 depicts the correlation between the S&P 500 Index and average CEO pay between 1970 and 2011. As shown in the figure, while "non-equity compensation" is at most weakly related to the performance of the overall stock market, total compensation was almost perfectly correlated until 2003, when the "bull market" from 2003-2007 was associated with relatively modest increases in average CEO pay.

We would expect *realized* compensation to vary with the overall market, since the gains from exercising non-indexed stock options will naturally increase with the market. But the compensation data in Figure 18 are based on the *grant-date* cost of the options, and not the amounts realized from exercising options. If compensation committees focused on the grant-date cost of options, we would expect the number of options granted to decrease when share prices increase, and would expect no systematic

---

[104] Based on a sample of approximately 600 large companies granting options to their CEOs during fiscal 1992, Murphy (1996) shows that about one-third of the companies reported grant-date values, while the remaining two-thirds reported potential values. Companies with higher dividend yields and lower volatilities (both factors that decrease Black–Scholes values) were significantly more likely to report Black–Scholes rather than potential values.

**Figure 18** Grant-Date Pay for CEOs in S&P 500 Firms vs. S&P 500 Index, 1970–2011. *Note:* The S&P 500 Index is a capitalization-weighted index of the prices of 500 large-cap common US stocks; the figure depicts monthly values. Compensation data are based on all CEOs included in the S&P 500, using data from Forbes and ExecuComp. CEO total pay includes cash pay, restricted stock, payouts from long-term pay programs and the value of stock options granted (using company fair-market valuations, when available, and otherwise using ExecuComp's modified Black–Scholes approach). Equity compensation prior to 1978 estimated based on option compensation in 73 large manufacturing firms (based on Murphy (1985)), equity compensation from 1979 through 1991 estimated as amounts *realized* from exercising stock options during the year, rather than grant-date values. Non-equity incentive pay is based on actual payouts rather than targets. Dollar amounts are converted to 2011-constant dollars using the Consumer Price Index.

correlation between the average pay and average market returns. However, if compensation committees focused on the number of options (e.g. granting the same number of options each year, as opposed to the same "value" of options each year), we would obtain the pattern in Figure 18. Because the grant–date Black–Scholes cost of an option is approximately proportional to the level of the stock price, awarding the same number of options after a doubling of stock prices amounts to doubling the cost of the option award. Therefore, if the number of options granted stayed constant over time, the cost of the annual option grants would have risen and fallen in proportion to the changes in stock prices.

If my interpretation of the data is correct, then the focus on the quantity (rather than cost) of options changed around 2002–2003. As I will argue below in Section 3.8.4, companies began voluntarily expensing the cost of options in 2002, both in response to the recent accounting scandals and in anticipation of mandated expensing in 2006.

In addition, in 2006 the SEC changed its disclosure rules to require option costs (rather than the number of options) in the Summary Compensation Table.

### 3.7.6  New York Stock Exchange Listing Requirements

Another contributing factor to the explosion in stock options—both to top executives and lower-level employees—was a 1998 change or "clarification" to New York Stock Exchange (NYSE) listing requirements. Under listing rules in affect at the time, companies needed shareholder approval for equity plans covering top-level executives, but did not need approval for broad-based plans. While the SEC had not been clear on how "broad-based" was defined, the general understanding was that such plans involved equity or option grants to employees below the executive level.

In January 1998, the NYSE quietly filed with the SEC a proposal clarifying definition of a "broad-based" plan as any plan in which (1) at least 20% of the company's employees were eligible to participate, and (2) at least half of the eligible employees were neither officers nor directors. The definition was a "safe harbor" (i.e. sufficient but not necessary): plans meeting the two criteria were presumed to be broadly based (and therefore could be introduced without shareholder approval), while plans falling outside these parameters would be considered on a case-by-case basis. The SEC received no letters questioning the proposed rule during the "public comment" period, and the ruling was approved and took effect on April 8, 1998. The final ruling was a surprise to shareholder advocates and institutions, who admitted to being embarrassed to have missed the proposal filing, and furious that it had been "buried" in the federal register and listed as a "cryptic notice" on the SEC's website.[105] Many observers speculated that the new rule was designed to lure NASDAQ companies to the NYSE, and most feared it would "open the floodgates" for executive stock options, since companies could avoid a shareholder vote by rolling their management plans into new broader-based plans. Consistent with my conclusions in Section 3.7.5, shareholder criticism focused exclusively on the dilutive effect of the option plans, on not on the transfer of value from shareholders to employees.

The NYSE—facing a barrage of criticism over its new rule—reopened the comment period (this time receiving 166) and created a task force to consider the new comments and make further suggestions. In June 1999, based on the recommendations of the task force, the NYSE issued "interim" new rules. Under the revised rules, the majority of the firm's non-exempt (e.g. non–managerial) employees (rather than 20% of all employees) must be eligible to participate, and the majority of options granted must go to non-officers (rather than the majority of the participants being non-officers). The new rule was an "exclusive test" rather than a safe harbor.

---

[105] Bryant, "New Rules on Stock Options By Big Board Irk Investors," *New York Times* (1998), Scipio, "NYSE Opens Option Loop Hole," *Investor Relations Business* (1998).

**Figure 19** Grant-Date Number of Employee Stock Options (measured by % of company shares) in the S&P 500, 1992–2005. *Note:* Figure shows the grant-date number of options as a percent of total common shares outstanding granted to all employees in an average S&P 500 firm, based on data from S&Ps ExecuComp database. Grants below the Top 5 executives are estimated based on Percent of Total Grant disclosures. Companies not granting options to any of their top five executives are excluded.

The new rules were enacted as companies faced growing political pressure to push grants to managers and employees at lower levels in the organization.[106] Several bills that encouraged broad-based stock option plans were introduced in Congress, including the Employee Stock Option Bill of 1997 (H.R. 2788) to ease the restrictions on qualified Incentive Stock Options granted to rank-and-file workers. At the same time, employees clamored for broad-based grants, but only if the company would promise that other components of their compensation would not be lowered. As a result of these pressures, the number and cost of options granted grew substantially.

Figure 19 shows the average annual option grants as a fraction of total common shares outstanding. In 1992, the average S&P 500 company granted its employees options on about 1.1% of its outstanding shares. In 2001, in spite of the bull market that increased share prices (that, in turn, increased the value of each granted option), the average S&P 500 company granted options to its executives and employees on 2.6% of its shares. By 2005, annual grants as a fraction of outstanding shares fell below 1995 levels to 1.3%.

---

[106] See, for example, Flanigan, "It's Time for All Employees to Get Stock Options," Los Angeles Times (1996), who argued that all employees should receive options if top executives receive options.

**Figure 20** Grant-Date Values of Employee Stock Options in the S&P 500, 1992–2005. *Note:* Figure shows the grant-date value of options (in millions of 2011-constant dollars) granted to all employees in an average S&P 500 firm, based on data from S&Ps ExecuComp data. Grants below the Top 5 are estimated based on Percent of Total Grant disclosures; companies not granting options to any of their top five executives are excluded. Grant-values are based on company fair-market valuations, when available, and otherwise use ExecuComps modified Black–Scholes approach.

Figure 20 shows the average inflation–adjusted grant–date values of options awarded by the average firm in the S&P 500 from 1992-2005.[107] Over this decade, the value of options granted increased from an average of $27 million per company in 1992 to nearly $300 million per company in 2000, falling to $88 million per company in 2005. Ignored in the news coverage and controversy over stock options awarded to CEOs and the next four highest-paid executives is the fact that employees and executives ranked below the top five have received between 85% and 90% of the total option awards.

Over the 14-year 1992–2005 time period, the average S&P 500 company awarded nearly $1.6 billion worth of options to its executives and employees (or nearly $800

---

[107] Options granted to lower-level executives and employees are estimated by dividing the options granted to the proxy-named executives by the percentage of all options that are granted to the proxy-named executives. Under the disclosure rules after 2006, the SEC no longer requires companies to report the percentage of all option awards that went to the proxy-named executives, and therefore my estimates of grants across the company end in 2005.

billion across all 500 companies). What is generally unappreciated is that in this process the average S&P 500 company transferred through options approximately 25.6% of its total outstanding equity to its executives and employees.[108]

Broad-based option grants were particularly generous in "new economy" firms and in firms below the S&P 500. Hall and Murphy (2003) show that the average new-economy firm in the S&P 500, S&P MidCap 400 and S&P SmallCap 600 granted options on 5.8% of its stock *annually* to employees below the top five between 1993 and 2001 (compared to only 2.3% annually in "old economy" firms). In 2000 alone, the average employee (below the top five) in the new-economy sector received options with a Black–Scholes value of $32000.[109]

The backlash against the explosion in option grants grew following the 2000 burst in the Internet bubble, when companies granted even more options at a lower price so that employees were not penalized for poor performance. Shareholder activists concerned about dilution pressured the NYSE to reconsider their rules. In late 2002, the NYSE and NASDAQ passed uniform new rules requiring shareholder approval for *all* equity plans, with no exemption for broad-based plans. The new rules—which also required shareholder approval for option repricings—were approved by the SEC and went into effect in July 2003.

## 3.8 The Accounting and Backdating Scandals (2001–2007)

### 3.8.1 Accounting Scandals and Sarbanes-Oxley

Accounting scandals erupted across corporate America during the early 2000s, destroying the reputations of once-proud firms such as Enron, WorldCom, Qwest, Global Crossing, HealthSouth, Cendant, Rite-Aid, Lucent, Xerox, Tyco International, Adelphia, Fannie Mae, Freddie Mac, and Arthur Andersen. In the midst of these scandals, Congress quickly passed the sweeping Sarbanes-Oxley Act in July 2002, setting or expanding standards for accounting firms, auditors, and boards of directors of publicly traded companies. The Act was primarily focused on accounting irregularities and not on compensation. However, Congress could not resist the temptation to use the new law to further regulate executive pay.

First, in direct response to the forgiveness of certain corporate loans given to executives at Tyco International, Section 402 of Sarbanes-Oxley prohibited all personal loans to executives and directors, regardless of whether such loans served a useful and

---

[108] The 25% calculation simply sums the annual percentages in Figure 19. This calculation overstates the transfer of equity to the extent that some options are forfeited or expire worthless, and understates the transfer of equity to the extent that the overall base of shares expands as options are exercised or as the company offers additional shares.

[109] The average grant value is determined by dividing the total value of grants in each industry (after excluding grants to the top five executives) by the total number of employees in the industry.

legitimate business purpose. For example, prior to Sarbanes-Oxley, companies would routinely offer loans to executives to buy company stock, often on a non-recourse basis so that the executive could fulfill the loan obligations by returning the purchased shares.[110] Similarly, companies attracting executives would routinely offer housing subsidies in the form of forgivable loans, a practice made unlawful under the new regulations.[111] Finally, Sarbanes-Oxley is viewed as prohibiting company-maintained cashless exercise programs for stock options, where an executive exercising options can use some of the shares acquired to finance both the exercise price and income taxes due upon exercise.[112]

Second, Section 304 of Sarbanes-Oxley requires CEOs and CFOs to reimburse the company for any bonus or equity-based compensation received, and any profits realized from selling shares, in the twelve months commencing with the filing of financial statements that are subsequently restated as a result of corporate misconduct. This "clawback" provision of Sarbanes-Oxley—which was subsequently extended in the TARP legislation and Dodd–Frank Financial Reform Act discussed below—was notable mostly for its ineffectiveness. Indeed, in spite of the wave of accounting restatements that led to the initial passage of Sarbanes-Oxley, the first individual clawback settlement under Section 304 did not occur until more than five years later, when UnitedHealth Group's former CEO William McGuire was forced to return $600 million in compensation.[113] The SEC became more aggressive in 2009, launching two clawback cases (CSK Auto and Diebold, Inc.) where the targeted executives were not accused of personal wrongdoing.[114]

---

[110] Indeed, it is easy to show that a traditional at-the-money stock option is equivalent to a non-recourse loan to purchase company stock at a zero interest rate with no down payment. Loans to purchase stock that carry a positive interest rate or require an executive down payment are less costly to grant than traditional options, and deliver better incentives by both forcing executives to invest some of their own money in the venture and only providing payouts when the stock price appreciates by at least the interest charged on the loan. It is unfortunate that Congress prohibited these types of plans.

[111] Offering housing subsidies in the form of loans that are forgiven with the passage of time is preferable to a lump-sum subsidy, since the company can avoid paying the full subsidy if the executive leaves the firm before the loan is repaid or fully forgiven.

[112] Technically, cashless exercise programs are implemented by offering the executive a short-term bridge loan to finance the purchase of the shares, followed by open-market transactions to sell some of the shares to repay the loan. Subsequent to Sarbanes-Oxley, executives exercising options have turned to conventional banks for bridge-loan financing, significantly increasing the transaction costs and further diluting the shares outstanding (since under company-maintained programs, the company need only issue the net number of shares and not the full number of shares under option).

[113] Plitch, "Paydirt: Sarbanes-Oxley A Pussycat On 'Clawbacks'," *Dow Jones Newswires* (2006), Bowe and White, "Record Payback over Options," *Financial Times* (2007).

[114] Berman, "The Game: New Frontier For the SEC: The Clawback," *Wall Street Journal* (2010), Korn, "Diebold to Pay $25 Million Penalty," *Wall Street Journal*(2010).

Finally, Section 403 of Sarbanes–Oxley required that executives disclose new grants of stock options within two business days of the grant; before the Act options were not disclosed until 10 days after the end of the month when the option was granted. As discussed in the next section, this provision had the unintended but ultimately beneficial effect of curbing option backdating for top executives more than two years before the existence of backdating was discovered.

### 3.8.2  Option Backdating

In 2005, academic research by University of Iowa professor Erik Lie and subsequent investigations by the *Wall Street Journal* unearthed a practice that became known as option backdating.[115] Under this practice, companies deliberately falsified stock option agreements so that options granted on one date were reported as if granted on an earlier date when the stock price was unusually low—commonly the lowest price in the quarter or in the year. Thus, options that were reported as granted *at* the money (that is, with an exercise price equal to the market price on the reported grant-date) were in reality granted *in* the money (that is, with an exercise price well below the market price on the actual grant-date). This unsavory practice violates federal disclosure rules, accounting and tax laws, and often violated the company's own stock-option policies, as follows:

- Under SEC rules in effect since 1993, companies granting options with an exercise price different from the fair market price on the grant-date are required to disclose this information to shareholders. Thus, companies backdating options should have informed shareholders that the options were actually issued with an exercise price less than the fair market value on the actual grant-date.
- As discussed in Sections 3.5.3 and 3.7.4, under FASB rules in effect before 2006, companies would typically face an accounting charge for stock options only if the exercise price was set lower than the grant-date market price. Thus, companies that backdated options reported no accounting expense when the actual accounting expense should have been the spread between the market and exercise price (amortized over the vesting period of the option). Companies backdating options are therefore not only falsifying option agreements, they are committing accounting fraud.
- As discussed in Section 3.7.3, compensation for proxy-named executives in excess of $1 million is deductible only if the compensation is performance based under the

---

[115] Key references include Lie (2005), Heron and Lie (2006b), Heron and Lie (2006a), Maremont, "Authorities probe improper backdating of options: Practice allows executives to bolster their stock gains; a highly beneficial pattern," *Wall Street Journal* (2005), Forelle and Bandler, "Backdating probe widens as two quit Silicon Valley firm; Power Integrations Officials leave amid options scandal; 10 companies involved so far," *Wall Street Journal* (2006), Forelle, "How Journal Found Options Pattern," *Wall Street Journal* (2006), "Hot Topic: Probing Stock-Options Backdating," *Wall Street Journal* (2006).

definition of IRS Section 162(m). In order for payments related to stock options to be considered performance based, the options must meet several criteria including having an exercise price that is at least as high as the grant–date market price.[116] Thus, assuming that the affected executives are subject to the $1 million threshold, companies that backdated options are taking deductions for compensation that is not deductible.

• Finally, most shareholder-approved stock option agreements include provisions specifying that option exercise prices must be no less than 100% of the market price on the date of grant. Thus, companies with such provisions that backdate options are violating their own internal policies.

The *Wall Street Journal's* crusade against backdating triggered SEC investigations into more than 140 firms. By August 2009, the SEC had filed civil charges against 24 companies and 66 individuals for backdating-related offenses, and at least 15 people had been convicted of criminal conduct.[117] In May 2007, Comverse Technology's former general counsel, William Sorin, pleaded guilty to a conspiracy charge and became the first corporate executive sent to prison for backdating executive and employee stock options; his boss (Comverse's founder and former CEO Kobi Alexander) fled to Namibia and is fighting extradition while remaining on the FBI's most wanted list.[118] In January 2008, Brocade's former CEO, Gregory Reyes, became the first executive to go to trial and be convicted on backdating charges; Reyes was sentenced to 21 months in prison and ordered to pay a $15 million fine. Brocade's former human resource executive was also convicted.[119] Reyes' conviction was thrown out by the US Court of Appeals in 2009, citing prosecutorial misconduct, but he was retried, reconvicted, and resentenced to 18 months in prison in June 2010.[120] In addition to the SEC civil and criminal charges, scores of companies have restated their financials based on internal investigations into backdating, and many have settled class action or derivative suits brought by shareholders.

Some backdating cases were obvious in retrospect, such as Cablevision's award of backdated options to its vice chairman after his death in 1999.[121] In most cases, however,

---

[116] If the amount of compensation the employee will receive under the grant or award is not based solely on an increase in the value of the stock after the date of grant or award (e.g. in the case of restricted stock, or an option that is granted with an exercise price that is less than the fair-market value of the stock as of the date of grant), none of the compensation attributable to the grant or award is qualified performance-based compensation. Internal Revenue Service, Section 1.162–27.

[117] Maremont, "Backdating Likely More Widespread," *Wall Street Journal*(2009).

[118] Bray, "Former Comverse Official Receives Prison Term in Options Case," *Wall Street Journal* (2007), "Fugitive Mogul's Rent Coup," *New York Post* (2009).

[119] Scheck and Stecklow, "Brocade Ex-CEO Gets 21 Months in Prison," *Wall Street Journal* (2008).

[120] Egelko, "18 months for ex-Brocade CEO," *San Francisco Chronicle* (2010).

[121] See Grant, Bandler and Forelle, "Cablevision Gave Backdated Grant To Dead Official," *Wall Street Journal* (2006).

executives would often go to considerable lengths to hide the backdating practices from the company's auditors, shareholders, and tax authorities. For example, in its investigation of backdating at Sycamore Networks, the SEC uncovered an internal menu that discussed ways to alter employees hire dates so they could get options with lower exercise prices, and also evaluated the risk that the changes might be discovered by auditors.[122] Executives at Mercury Interactive used WhiteOut to alter the dates on option documents, and joked about magic backdating ink.[123]

As noted above in Section 3.8.1, changes in reporting requirements in 2002 essentially put an end to option backdating for top-level executives more than two years before academics and the media uncovered the practice. Between May 1992 and August 2002, option grants for corporate insiders were typically not disclosed until 10 days after the end of the month when the option was granted, providing substantial opportunity for manipulating grant-dates. In August 2002, as part of the Sarbanes-Oxley Act, the SEC required executives receiving options to disclose those grants within two business days after the grant was made. Heron and Lie (2006a) and Narayanan and Seyhun (2005) show that the abnormal run-up in stock prices following reported grant-dates (which they interpret as evidence of backdating) declined substantially after the new reporting rules, thus suggesting that the Sarbanes-Oxley Act had the unintended (but desirable) effect of stemming backdating practices.[124]

By 2010, the SEC's investigations and prosecutions of backdating had wound down. New disclosure rules introduced in 2006 were designed to identify new backdating cases by requiring companies to report not only exercise prices for option grants, but also the grant-date market price, date of grant, and the date that the board approved the grant.[125] While there is no accepted count of the number of companies engaged in backdating (beyond the 24 companies formally charged by the SEC or the approximately 150 companies that have restated financials after internal investigations revealed backdating[126]), academic research has suggested that the practice was widespread. Based on statistical analysis of exercise prices, Edelson and Whisenant (2009) estimate that as many as 800 firms engaged in the practice; other estimates have been as high as 2000.[127]

---

[122] Hechinger and Bandier, "In Sycamore Suit, Memo Points to Backdating Claims," *The Wall Street Journal* (2006). The internal memo is available at Sycamore Networks (2001).

[123] See Lee, "Option lawsuit give up details: Shareholders suing Mercury Interactive over timing of grants," *San Francisco Chronicle* (2007).

[124] The reporting requirements under Sarbanes-Oxley apply only to executive officers and directors, and there is evidence from SEC investigations that some companies continued backdating for lower-level employees subsequent to the August 2002. However, since grants to such employees are not publicly disclosed, it has not been possible to perform a comprehensive analysis of the practice.

[125] In the proxy disclosure rules in effect between 1993 and 2006, companies were required to report the expiration date for new grants, but not the grant date.

[126] Nicklaus, "Scandal left both sides sullied: Backdating undermined confidence, but some 'good guys' overreached," *St. Louis Post-Dispatch* (2010).

[127] Ryst, "How To Clean Up A Scandal," *BusinessWeek.com* (2006).

In retrospect, while issuing options with exercise prices below grant–date market prices can be part of an efficient compensation structure, it is difficult to defend the practice of backdating and the ex post manipulation and falsification of grant–dates. However, it is also difficult to defend the SEC's aggressiveness in prosecuting and criminalizing what would seem to be relatively minor books and records infractions. Consider the following:

- There is nothing illegal about setting exercise prices to the lowest price observed during a month or quarter (or any other price), as long as the company appropriately discloses the practice and (based on FASB rules in effect before 2006) records an accounting expense equal to the difference between the exercise price and the market price on the true grant–date. In practice, however, very few firms issue options with exercise prices below market prices precisely because of the accounting charge associated with such options.

- Companies charged with backdating have restated their financials to reflect the actual spread between the exercise and market price. However, this remedy misses the point: the relevant alternative to backdating was *not* issuing in–the–money options and taking an accounting charge, but rather issuing a larger number of at–the–money options and avoiding the accounting charge. Therefore, under this relevant alternative, there would be no change in reported accounting expenses or earnings, but there would be an increase in the number of options granted.

- There is no evidence to my knowledge that companies engaged in backdating systematically overpaid lower–level employees receiving such grants, thus no evidence that backdating was associated with a large transfer of wealth from shareholders to employees.[128]

The SEC prosecuted backdating cases with a zeal usually reserved for hardened criminals. Executives associated with backdating schemes were charged with myriad crimes, including filing false documents, securities fraud, and conspiracy to commit securities fraud. KB Homes former CEO Bruce Karatz, for example, faced up to 415 years in prison if convicted on all backdating-related charges including 15 counts of mail, wire, and securities fraud, four counts of making false statements in SEC filings, and one count of lying to his company's accountants. Mr. Karatz was ultimately convicted in April 2010 on two counts of mail fraud, one count each of making false statements in SEC filings and to his accountants and faced up to 80 years in prison.[129] Ultimately, however, Mr. Karatz was sentenced to five years probation (including eight months of house arrest), a $1 million fine and 2,000 hours of community service.

---

[128] Bebchuk, Grinstein, and Peyer (2010) show that CEOs receiving lucky grants (which they define as grants with exercise prices set at the lowest price during the grant month) have higher total compensation than CEOs without lucky grants.

[129] Wotapka, "Former CEO At KB Home Is Convicted," *Wall Street Journal* (2010).

The SEC's record of successful convictions has been far from perfect. Its suit against Michael Shanahan for backdating at Engineered Support Systems was dismissed mid-trial when the judge determined that the SEC's case provided no evidence of fraud. Similarly, the SEC's high-profile case against Broadcom was dismissed amid claims of significant prosecutorial misconduct and lack of criminal intent.[130]

### 3.8.3  Enron and Section 409(A)

Enron, like many other large companies, allowed mid-level and senior executives to defer portions of their salaries and bonuses through the company's non-qualified deferred compensation program. When Enron filed for Chapter 11 bankruptcy protection in December 2002, about 400 senior and former executives became unsecured creditors of the corporation, eventually losing most (if not all) of the money in their accounts.[131] However, just before the bankruptcy filing, Enron allowed a small number of executives to withdraw millions of dollars from their deferred compensation accounts. The disclosure of these payments generated significant outrage (and law suits) from Enron employees who lost their money, and attracted the ire of Congress.

As a direct response to the Enron situation, Section 409(A) was added to the Internal Revenue Code as part of the American Jobs Creation Act of 2004. In essence, the objectives of Section 409(A) were to limit the flexibility in the timing of elections to defer compensation in nonqualified deferred compensation programs, to restrict withdrawals from the deferred accounts to pre-determined dates (and to prohibit the acceleration of withdrawals), and to prevent executives from receiving severance-related deferred compensation until six months after severance. Section 409(A) imposes taxes on individuals with deferred compensation as soon as the amounts payable under the plan are no longer subject to a substantial risk of forfeiture. Individuals failing to pay taxes in the year the amounts are deemed to no longer be subject to the substantial forfeiture risk owe a 20% excise tax and interest penalties on the amount payable (even if the individual has not received or may never receive any of the income).

One of the notable features of Section 409(A) is that it significantly broadens the definition of deferred compensation. For example, annual bonuses or reimbursement of expenses paid more than two and a half months after the close of the fiscal year are considered deferred compensation subject to Section 409(A). Similarly, supplemental executive retirement plans (SERPs), phantom stock awards, stock appreciation rights, split-dollar life insurance arrangements, and individual employment agreements

---

[130] Henning, "Behind the Fade-Out of Options Backdating Cases," *New York Times* (2010).

[131] Barboza, "Enron's Many Strands: Executive Compensation. Enron paid some, not all, deferred compensation," *New York Times* (2002).

allowing deferral of compensation or severance awards are also (under some circum-stances) considered deferred compensation subject to Section 409(A).

While developed as a response to the Enron situation, Section 409(A) was still being drafted when the option backdating scandals came to light. As a result, Congress defined discount options (i.e. options with an exercise price below the market price on the date of grant) as deferred compensation subject to Section 409(A). In particular, Section 409(A) requires discount options to have a fixed exercise date (that is, a date in the future when the option must be exercised). Unless the option holder pre-commits to the future date when the option will be exercised, the holder is subject to a 20% penalty tax, in additional to regular income tax, plus possible interest and other penalties, regard-less of whether the option is ever exercised.[132] The new rule applied retroactively to options granted before 2005 but not vested as of December 31, 2004, and was explicitly designed to penalize senior executives receiving backdated options.

### 3.8.4  Accounting for Options (Finally!) and the Rise of Restricted Stock

The first decade of the new century have brought several important changes in the level and composition of CEO pay. As shown in Figure 21, median grant-date total CEO pay in the S&P 500 declined from $9.3 million in the peak year of 2001 to $9.0 million in 2011, representing the first prolonged stagnation in CEO pay since the early 1970s. The decrease in pay primarily reflects both a substantial decline in the grant-date value of stock options, and a shift in the industry composition of the S&P 500. In 2001, the value of stock options at the award date accounted for 53 percent of the pay for the typical S&P 500 CEO. By 2011, options accounted for only 21 percent of total pay. Moreover, the decline in stock option grants in the early 2000s has been associated with an increase in stock grants, which accounted for 36% of average pay by 2011 (up from only 8% in 2001). The stock grants include a mixture of traditional restricted stock (vesting only with the passage of time) and performance shares (where vesting is based on performance criteria).

Figure 22 shows the percentage of S&P 500 companies that made stock option or restricted stock grants to their CEOs between 1992 and 2011. The percentage of compa-nies granting options to their CEOs in each year increased from about 63% in 1992 to 87% by 2001, falling to 68% in 2011. Over the same time period, the percentage of companies

---

[132] IRS guidance has not been clear with respect to the amount subject to the additional 20% penalty. For example, Morrison and Foerster (http://www.mofo.com/news/updates/files/update02204.html) has advised its clients that the amount subject to the penalty could be any of the following: the difference between the exercise price and the fair market value of the stock subject to the option measured on the date of grant of the option; the difference between the exercise price and the fair market value of the stock subject to the option measured on the date the shares subject to the option vest; the difference between the exercise price and the fair market value of the stock subject to the option measured on the date of exercise; the Black-Scholes value of the option measured on the date of grant of the option; or the Black-Scholes value of the option measured on the date the shares subject to the option vest.

**Figure 21** Median Grant-date Compensation for CEOs in S&P 500 Firms, 2001–2011. *Note:* Median pay levels (in 2011-constant dollars) based on ExecuComp data for S&P 500 CEOs. Total compensation (indicated by bar height) defined as the sum of salaries, non-equity incentives (including bonuses), benefits, stock options (valued on the date of grant using company fair-market valuations, when available, and otherwise using ExecuComps modified Black–Scholes approach), stock grants, and other compensation. Other compensation excludes pension-related expenses. Pay composition percentages defined as the average composition across executives.

making restricted stock or performance–share grants more than tripled from 25% to 82%. The trend suggests a substitution of stock grants for stock options, although more than half of the S&P 500 CEOs have received *both* options and restricted stock annually since 2006.

One obvious explanation for the drop in stock options and the rise in restricted stock since the early 2000s is the stock market crash associated with the burst of the Internet Bubble in 2000 and exacerbated by the terrorist attacks on the World Trade Center in 2001. In particular, the sharp market–wide decline in stock prices in the early 2000s left many outstanding options underwater and lowered executive expectations for the future increases in their company's stock prices. Indeed, in many cases, including Microsoft and Cablevision, current outstanding (but out–of–the–money) options were cancelled and replaced with restricted stock, often at terms very favorable to executives. Executives will naturally prefer restricted stock to options when they have low expectations for future firm performance. While restricted stock will always retain some value as long as the firm is valued at greater than its liabilities, executives often expect that options granted in a declining market are likely to expire worthless.

Indeed, stock options have always become more popular when stock markets are trending upward (i.e. bull markets) and less popular when markets trend down

**Figure 22** CEOs in S&P 500 Firms receiving equity-based compensation, 1992–2011. *Note:* Sample is based on all CEOs included in the S&P 500, based on S&P's ExecuComp database. Stock grants include both restricted and performance shares.

(i.e. bear markets). As documented throughout this history of CEO pay, almost every recession over the past 60 years has been associated with a reduced use of stock options, and during the lackluster 1970s many firms replaced their option plans with new accounting-based bonus plans designed to provide more predictable payouts. However, the spike in the importance of restricted shares in 2006 (rising in Figure 21 from 17% to 26% of total pay from 2005) in a year with robust stock-market performance (the Dow Jones increased by 16% in 2006) suggests that the decline in stock options in favor of restricted shares reflects more than market trends. I believe the answer largely reflects changes in the accounting treatment of options.

The scandals that erupted across corporate America during the early 2000s focused attention on the quality of accounting disclosures, which in turn renewed pressures for companies to report the expense associated with stock options on their account–ing statements. Before 2002, only a handful of companies had elected to expense options under FAS123; the remainder elected to account for options under the old rules (where there was typically no expense). In the summer of 2002, several dozen firms announced their intention to expense options voluntarily; more than 150 firms had elected to expense options by early 2003 (Aboody, Barth, and Kasznik, 2004). Moreover, shareholder groups (most often representing union pension funds) began demanding shareholder votes on whether options should be expensed; more than

150 shareholder proposals on option expensing were submitted during the 2003 and 2004 proxy season (Ferri and Sandino, 2009). By late 2004, about 750 companies had voluntarily adopted or announced their intention to expense options. In December 2004, FASB announced FAS123R which revised FAS123 by *requiring* all US firms to recognize an accounting expense when granting stock options, effective for fiscal years beginning after June 15, 2005.

In addition to requiring an accounting expense for all options granted after June 15, 2005, FAS123R required firms to record an expense for options granted *before* this date that were not yet vested (or exercisable) as of this date. To avoid taking an accounting charge for these outstanding options, many firms accelerated vesting of existing options so that all options were exercisable by June 15, 2005 (Choudhary et al., 2009).

Under the accounting rules in place since 1972 (and continuing under FAS123R), companies granting traditional restricted stock (vesting only with the passage of time) recognize an accounting expense equal to the grant-date value of the shares amortized over the vesting period. Under FAS123R, the expense for stock options is similar to that of shares of stock: companies must recognize an accounting expense equal to the grant-date value of the options amortized over the period when the option is not exercisable. Option expensing (whether voluntarily under FAS123, or by law under FAS123R) significantly leveled the playing field between stock and options from an accounting perspective. As a result, companies reduced the number of options granted to top executives (and other employees), and greatly expanded the use of restricted shares.

The new accounting rules also facilitated another change long desired by shareholder advocates: a switch from traditional time-lapse restricted stock to "performance shares" that vest only upon achievement of accounting- or market-based performance goals. Angelis and Grinstein (2011), for example, report that 52% of the 2007 restricted stock awards for CEOs in the S&P 500 were performance-based.

Under the 1972 rules, performance shares were expensed using "variable" rather than "fixed" accounting, meaning that the company would record an expense based on the grant-date stock price, and then record additional expenses reflecting the appreciation or depreciation of the performance share up until the date that the performance hurdle was achieved. Therefore, if the stock price increased between the grant and the achievement of the performance hurdle (which is typically the case), the accounting expense for performance shares was higher than the accounting expense for traditional time-lapse restricted stock. In contrast, under FAS123R fair-market-value accounting, the expense for performance shares is generally less than the expense for traditional restricted stock, because the company can take into account the severity of the performance hurdles when estimating the fair-market value. In addition, while traditional restricted stock is considered non-performance-related under IRS Section 162(m) (and

thus subject to the $1 million deductibility cap), performance shares can be structured to be fully deductible.

### 3.8.5  Conflicted Consultants and CEO Pay[133]

Most large companies rely on executive compensation consultants to make recommendations on appropriate pay levels, to design and implement short-term and long-term incentive arrangements, and to provide survey and competitive-benchmarking information on industry and market pay practices. In addition, consultants are routinely asked to opine on existing compensation arrangements and to give general guidance on change-in-control and employment agreements, as well as on complex and evolving accounting, tax, and regulatory issues related to executive pay.

Critics seeking explanations for high executive pay have increasingly accused these consultants as being (partly) to blame for the perceived excesses in pay. Concerns over the role of consultants led the SEC – as part of their 2006 overhaul of proxy disclosure rules – to require companies to identify any consultants who provided advice on executive or director compensation; to indicate whether or not the consultants are appointed by the companies' compensation committees; and to describe the nature of the assignments for which the consultants are engaged.

Initial results from the 2007 proxy season appeared to buttress the concerns of the critics. An October 2007 report issued by the Corporate Library, "The Effect of Compensation Consultants" (Higgins, 2007) concluded that companies using consultants offer significantly higher pay than companies not using consultants.[134] However, the cross-sectional correlation between CEO pay and the use of consultants does not imply that the consultants *caused* the high pay; it is equally plausible that companies with high pay are most likely to seek the advice of consultants. Indeed, Armstrong, Ittner, and Larcker (2012) find no evidence of differences in pay between a sample of firms using consultants and a matched sample of firms not using consultants. Similarly, based on a time-series of 2006-2009 data, Murphy and Sandino (2012) find no evidence that firms increase pay after retaining consultants.

The SEC's disclosure requirements were followed by Congressional hearings and a December 2007 report from the US House of Representatives Committee on Oversight and Government Reform, "Executive Pay: Conflicts of Interest Among Compensation Consultants" (Waxman, 2007). The Congressional hearings focused on consultants offering a full range of compensation, benefits, actuarial and other human resources services in addition to executive pay. The provision of these other services creates a potential conflict of interest because the decisions to engage the consulting firm

---

[133] This section draws heavily from Murphy and Sandino (2010, 2012).

[134] Academic studies based on the first year of consultant disclosures – including Cadman, Carter, and Hillegeist (2010), Armstrong, Ittner, and Larcker (2012) and (early versions of) Murphy and Sandino (2010) – also documented significantly higher pay in companies using consultants.

in these more-lucrative corporate-wide consulting areas are often made or influenced by the same top executives who are benefited or harmed by the consultant's executive pay recommendations.

In response to the Congressional concerns, the SEC expanded its disclosure rules in 2009 to require firms to disclose fees paid to their executive compensation consultants whenever the consultants received more than $120,000 for providing any other services to the firm beyond those related to executive and director pay. The SEC exempted from these requirements firms that retain at least one compensation consultant that works exclusively for the board, and also exempted disclosing consultants that affect executives' and directors' compensation only through providing advice related to broad-based plans that do not discriminate executives and/or directors from other employees. As discussed below in Section 3.10.2, the SEC disclosure rules were further expanded in 2012 (as part of the implementation of the Dodd–Frank Act) to require firms to disclose whether the work of the consultant has raised any conflict if interest and, if so, the nature of the conflict and how the conflict is being addressed.

The initial and expanded SEC disclosure rules were introduced without any evidence that "conflicted consultants" were, indeed, complicit in perceived pay excesses. Based on the initial year of consultant disclosures, Cadman, Carter, and Hillegeist (2010) find no evidence that CEO pay is related to consultant conflicts of interest. Based on similar data (supplemented with IRS and Department of Labor data identifying actuarial service providers), Murphy and Sandino (2010) find some evidence that CEO pay is modestly higher in firms where consultants provide other services. However, in subsequent time-series analyses, Murphy and Sandino (2012) show that the relation between conflicted consultants and CEO pay had become statistically and economically insignificant by 2008.

While the evidence suggests, at most, a modest link between conflicted consultants and CEO pay, the SEC disclosure requirements have resulted in dramatic changes in the compensation consulting industry. The largest full-service consulting firms in 2006 (Towers Perrin, Mercer, Hewitt, and Watson Wyatt) have experienced significant declines in market share among their S&P 500 clients, while the largest non-integrated firms focused only on executive compensation (Frederick Cook and Co. and Pearl Meyer) have increased market share. In addition, many of the top consultants from the full-service firms left to create their own "boutique" firms focused on advising boards. For example, consultants from Towers Perrin and Watson Wyatt formed Pay Governance, consultants from Hewitt formed Meridian Compensation Partners, and consultants from Mercer formed Compensation Advisory Partners. The full-service firms have also consolidated: Towers Perrin and Watson Wyatt merged to create Towers Watson, while Hewitt was acquired by Aon.

As discussed by Murphy and Sandino (2010), the experience of the full-service consulting firms closely parallels the experience of accounting firms offering both auditing and consulting services. Concerns regarding conflicts when accounting firms

offered services beyond auditing led not only to the Sarbanes-Oxley Act and to detailed disclosures of fees charged for auditing and non-auditing businesses, but also to the practice of companies avoiding using their auditors for other services. This practice has defined the industry, in spite of the fact that the auditors (with their vast firm-specific knowledge) might be the efficient provider of such services, and notwithstanding the fact that there was no direct evidence that these potential conflicts actually translated into misleading audits.

## 3.9  Pay Restrictions for TARP Recipients (2008–2009)

### 3.9.1  The Emergency Economic Stabilization Act (EESA)

On September 19, 2008—at the end of a tumultuous week on Wall Street that included the Lehman Brothers bankruptcy and the hastily arranged marriage of Bank of America and Merrill Lynch—Treasury Secretary Paulson asked Congress to approve the Administration's plan to use taxpayers' money to purchase "hundreds of billions" in illiquid assets from US financial institutions.[135] Paulson's proposal contained no constraints on executive compensation, fearing that restrictions would discourage firms from selling potentially valuable assets to the government at relatively bargain prices.[136] Limiting executive pay, however, was a long-time top priority for Democrats and some Republican congressmen, who viewed the "Wall Street bonus culture" as a root cause of the financial crisis. Congress rejected the bailout bill on September 30, but reconsidered three days later after a record one-day point loss in the Dow Jones Industrial Average and strong bipartisan Senate support. The Emergency Economic Stabilization Act (EESA) was passed by Congress on October 3[rd], and signed into law by President Bush on the same day.

When Treasury invited (or, in some cases, coerced) the first eight banks to participate in TARP, a critical hurdle involved getting the CEOs and other top executives to waive their rights under their existing compensation plans. At the time, the proposed restrictions seemed serious. For example, while Section 304 of the 2002 Sarbanes-Oxley Act required clawbacks of certain executive ill-gotten incentive payments, the Act only covered the CEO and chief financial officer (CFO), and only covered accounting restatements. While applying only to TARP recipients (Sarbanes-Oxley applied to all firms), the October 2008 EESA covered the top-five executives (and not just the CEO and CFO), and covered a much broader set of material inaccuracies in performance metrics. In addition, EESA lowered the IRS cap on deductibility for the top-five executives from $1 million to $500,000, and applied this limit to all forms of compensation (and not

[135] Solomon and Paletta, "US Bailout Plan Calms Markets, But Struggle Looms Over Details", *Wall Street Journal* (2008).

[136] Hulse and Herszenhorn, "Bailout Plan Is Set; House Braces for Tough Vote", *New York Times*(2008).

just non–performance-based pay). EESA also prohibited new golden parachutes agreements for the Top 5 executives, and capped payments under existing plans to 300% of the executives' average taxable compensation over the prior five years.

In a semantic change that will confuse students of executive compensation for years to come, EESA also formally defined "golden parachutes" as amounts paid in "the event of an involuntary termination, bankruptcy filing, insolvency, or receivership". Previously, the term "golden parachute" had referred exclusively (if not pejoratively) to payments made in connection with a change in control. Under IRS Section 280(G) (discussed above in Section 3.6.1), change–in–control payments exceeding 300% of executives' average taxable compensation over the prior five years were subject to significant tax penalties. Thus, EESA not only explicitly capped payments, but substantially expanded the events characterized as golden parachutes.

### 3.9.2  The American Reinvestment and Recovery Act (ARRA) Amends EESA

In January 2009, reports began surfacing that Merrill Lynch distributed $3.6 billion in bonuses to its 36,000 employees just before the completion of the merger with Bank of America: the top 14 bonus recipients received a combined $250 million, while the top 149 received $858 million (Cuomo, 2009). The CEOs of Bank of America and the former Merrill Lynch (neither of whom received a bonus for 2008) were quickly hauled before Congressional panels outraged by the payments, and the Attorney General of New York launched an investigation to determine if shareholders voting on the merger were misled about both the bonuses and Merrill's true financial condition. The SEC joined in with its own civil complaint, which sued the Bank of America but not its individual executives, and the bank agreed to settle for $33 million. However, a few weeks later a federal judge threw out the proposed settlement, insisting that individual executives be charged and claiming that the settlement did not comport with the most elementary notions of justice and morality.[137] In February 2010, the judge relented and reluctantly approved the settlement after it had been increased to $150 million.[138]

By the time the Merrill Lynch bonuses were revealed, the country had a new President, a new Congress, and new political resolve to punish the executives in the companies perceived to be responsible for the global meltdown. Indicative of the mood in Washington, Senator McCaskill (D-Missouri) introduced a bill in January 2009 that would limit total compensation for executives at bailed-out firms to $400,000, calling Wall Street executives "a bunch of idiots who were kicking sand in the face of the American taxpayer".[139]

---

[137] Scannell, Rappaport, and Bravin, "Judge Tosses Out Bonus Deal—SEC Pact With BofA Over Merrill Is Slammed; New York Weighs Charges Against Lewis," *Wall Street Journal* (2009).

[138] Fitzpatrick, Scannell, and Bray, "Rakoff Backs BofA Accord, Unhappily," *Wall Street Journal* (2010).

[139] Andrews and Bajaj, "Amid Fury, US Is Set to Curb Executives' Pay After Bailouts," *New York Times* (2009).

On February 4, 2009, President Obama's administration responded with its own proposal for executive-pay restrictions that distinguished between failing firms requiring exceptional assistance and relatively healthy firms participating in TARPs Capital Purchase Program. Most importantly, the Obama Proposal for exceptional assistance firms (which specifically identified AIG, Bank of America, and Citigroup) capped annual compensation for senior executives to $500,000, except for restricted stock awards (which were not limited, but could not be sold until the government was repaid in full, with interest). In addition, for exceptional-assistance firms the number of executives subject to clawback provisions would be increased from 5 under EESA to 20, and the number of executives with prohibited golden parachutes would be increased from 5 to 10. In addition, the next 25 highest-paid executives would be prohibited from parachute payments that exceed one year's compensation.

Moreover—in response to reports of office renovations at Merrill Lynch, corporate jet orders by Citigroup, and corporate retreats by AIG—the Obama Proposal stipulated that all TARP recipients adopt formal policies on luxury expenditures. Finally, the Obama Proposal required all TARP recipients to fully disclose their compensation policies and allow nonbinding Say-on-Pay shareholder resolutions.[140]

In mid-February 2009, separate bills proposing amendments to EESA had been passed by both the House and Senate, and it was up to a small conference committee to propose a compromise set of amendments that could be passed in both chambers. On February 13—as a last-minute addition to the amendments—the conference chairman (Senator Chris Dodd) inserted a new section imposing restrictions on executive compensation that were opposed by the Obama administration and severe relative to both the limitations in the October 2008 version and the February 2009 Obama Proposal. Nonetheless, the compromise was quickly passed in both chambers with little debate and signed into law as the American Recovery and Reinvestment Act of 2009 by President Obama on February 17, 2009.

Table 1 compares the pay restrictions under the original 2008 EESA bill, the 2009 Obama Proposal, and the 2009 ARRA (which amended Section 111 of the 2008 EESA). While the clawback provisions under the original ESSA covered only the top-five executives (up from only two in SOX), the Dodd Amendments extended these provisions to 25 executives and applied them retroactively.[141] In addition, while the

[140] TARP recipients not considered exceptional assistance firms could waive the disclosure and Say-on-Pay requirements, but would then be subject to the $500,000 limit on compensation (excluding restricted stock).

[141] The number of executives covered by the Dodd Amendments varied by the size of the TARP bailout, with the maximum number effective for TARP investments exceeding $500 million. As a point of reference, the average TARP firm among the original eight recipients received an average of $20 *billion* in funding, and virtually all the outrage over banking bonuses have involved banks taking well over $500 million in government funds. Therefore, I report results assuming that firms are in the top group of recipients.

**Table 1**  Comparison of Pay Restrictions in EESA (2008), Obama Proposal (2009), and ARRA (2009)

**A.** *Limits on Pay Levels and Deductibility*

| | |
|---|---|
| Pre–EESA (IRS §162(m) (1994)) | Limits deductibility of top–5 executive pay to $1000,000, with exceptions for performance–based pay. |
| EESA (2008) All TARP Recipients | Limits deductibility of top–5 executive pay to $500,000, with no exceptions for performance-based pay. |
| Obama (2009) Exceptional Assistance Firms | In addition to deductibility limits, cash pay is capped at $500,000; additional amounts can be paid in restricted shares vesting after government paid back. |
| Obama (2009) Other TARP Recipients | Same as exceptional assistance firms, but pay caps can be waived if firm offers full disclosure of pay policies and a non–binding Say–on–Pay vote. |
| ARRA (2009) All TARP Recipients | In addition to deductibility limits, disallows all incentive payments, except for restricted stock capped at no more than one–half base salary. No caps on salary. |

**B.** *Golden Parachutes*

| | |
|---|---|
| Pre–EESA (IRS §280G (1986)) | Tax penalties for change–in–control–related payments exceeding 3 times base pay. |
| EESA (2008) Auction Program | No new severance agreements for Top 5. |
| EESA (2008) Capital Purchase Program | No new severance agreements for Top 5, and no payments for top 5 executives under existing plans exceeding 3 times base pay. |
| Obama (2009) Exceptional Assistance Firms | No payments for Top 10; next 25 limited to 1 times base pay. |
| Obama (2009) Other TARP Recipients | No payments for top 5 executives under existing plans exceeding 1 times base pay. |
| ARRA (2009) All TARP Recipients | No payments for Top 10. Disallows all payments (not just excess payments). |

**C.** *Clawbacks*

| | |
|---|---|
| Pre–EESA (Sarbanes–Oxley (2002)) | Covers CEO and CFO of publicly traded firms following restatements |
| EESA (2008) Auction Program | No new provisions. |
| EESA (2008) Capital Purchase Program | Top 5 executives, applies to public and private firms, not exclusively triggered by restatement, no limits on recovery period, covers broad material inaccuracies (not just accounting restatements). |
| Obama (2009) All TARP Recipients | Same as above, but covers 20 executives. |
| ARRA (2009) All TARP Recipients | Covers 25 executives for all TARP participants, retroactively. |

original ESSA disallowed severance payments in excess of 300% of base pay for the top five executives, the Dodd Amendments covered the top 10 executives and disallowed *all* payments (not just those exceeding 300% of base). Most importantly, the Dodd Amendments allowed only two types of compensation: base salaries (which were not restricted in magnitude), and restricted stock (limited to grant-date values no more than half of base salaries). The forms of compensation explicitly prohibited under the Dodd amendments for TARP recipients include performance-based bonuses, retention bonuses, signing bonuses, severance pay, and all forms of stock options.

Finally, the Dodd amendments imposed mandatory Say-on-Pay resolutions for all TARP recipients. In early 2009—not long after the Dow Jones Industrial Average hit its crisis minimum at about 6500—shareholders had an opportunity to provide a non-binding vote of approval on the 2008 compensation received by the top executives at the TARP recipients (i.e. compensation for the year when these firms allegedly dragged the economy into a financial crisis). As an interesting historical footnote, none of the TARP recipients received a majority vote against its executive compensation levels and policies.

As another interesting historical footnote: while almost all attempts to regulate executive compensation have produced negative unintended side affects, the Dodd Amendments produced a positive one. In particular, many TARP recipients found the draconian pay restrictions sufficiently onerous that they hurried to pay back the government in time for year-end bonuses.

As draconian as the Dodd Amendments (triggered by the Merrill Lynch payments) were, things were about to get worse. The second flash point for outrage over bonuses involved insurance giant American International Group (AIG), which had received over $170 billion in government bailout funds, in large part to offset over $40 billion in credit default-swap losses from its Financial Products unit. In March 2009, AIG reported it was about to pay $168 million as the second installment of $450 million in contractually obligated retention bonuses to employees in the troubled unit. (The public outrage intensified after revelations that most of AIGs bailout money had gone directly to its trading partners, including Goldman Sachs ($13 billion), Germanys Deutsche Bank ($12 billion), and France's Société Générale ($12 billion)). The political fallout was swift and furious: in the week following the revelations seven bills were introduced in the House and Senate aimed specifically at bonuses paid by AIG and other firms bailed out through TARP:

- H.R. 1518, the Bailout Bonus Tax Bracket Act of 2009 imposed a 100% tax on bonuses over $100,000.
- H.R. 1527 imposed an additional 60% tax (on top of 35% ordinary income tax) on bonuses exceeding $100,000 paid to employees of businesses in which the federal government has an ownership interest of 79% or more. (Not coincidentally, the government owned 80% of AIG when the bill was introduced).

- H.R. 1575, the End Government Reimbursement of Excessive Executive Disbursements Act (i.e. the End GREED Act) authorized the Attorney General to seek recovery of and limit excessive compensation.
- H.R. 1577, the AIG Bonus Payment Bill required the Secretary of Treasury to implement a plan within two weeks to thwart the payment of the AIG bonuses, and required Treasury approval of any future bonuses by any TARP recipient.
- H.R. 1586 sought to impose a 90% income tax on bonuses paid by TARP recipients; employees would be exempt from the tax if they returned the bonus in the year received.
- S. 651, the Compensation Fairness Act of 2009, imposed a 70% excise tax (half paid by the employee and half by the employer) for any bonus over $50,000 paid by a TARP firm.
- H.R. 1664, the Pay for Performance Act of 2009 prohibited any compensation payment (under existing as well as new plans) if such compensation: (1) is deemed unreasonable or excessive by the Secretary of the Treasury; and (2) includes bonuses or retention payments not directly based on approved performance measures. The bill also created a Commission on Executive Compensation to study and report to the President and Congress on the compensation arrangements at TARP firms.

Most of these bills were either stalled in committees or failed in a vote, although many features of H.R. 1664 were incorporated into the July 2010 Dodd–Frank Wall Street Reform bill discussed below. Therefore, the reason to list the bills above is not for their ultimate relevance to policy, but rather as evidence of Congressional outrage and a political resolve to punish Wall Street for its bonus practices.

While details on the compensation of the five highest-paid executive officers are publicly disclosed and widely available, banks have historically been highly secretive about the magnitude and distribution of bonuses for its traders and investment bankers. Indeed, since the SEC disclosure rules only apply to *executive officers*, the banks can have non-officer employees making significantly more than the highest-paid officers. Following the Merrill Lynch and AIG revelations, New York Attorney General Andrew Cuomo subpoenaed bonus records from the nine original TARP recipients, arguing that New York law allows creditors to challenge any payment by a company if the company did not get adequate value in return. His report—published in late July 2009—was provocatively titled: No Rhyme or Reason: The Heads I Win, Tails You Lose Bank Bonus Culture.

Table 2 summarizes the distribution of bonuses for the nine original TARP recipients, based on data from the Cuomo (2009) report. The table shows, for example, that 738 Citigroup employees received bonuses over $1 million, and 124 received over $3 million, in a year when the bank lost nearly $30 billion. The 2008 bonus pools exceeded annual earnings in six of the nine banks; in aggregate the banks paid $32.6 billion in

**Table 2** 2008 Earnings and Bonus Pools for Eight Original TARP Recipients

| Corporation | 2008 Earnings/ (Losses) ($bil) | 2008 Bonus Pool ($bil) | Number of Employees Receiving Bonuses Exceeding | | |
| --- | --- | --- | --- | --- | --- |
| | | | $3 mil | $2 mil | $1 mil |
| Bank of America | $4.0 | $3.3 | 28 | 65 | 172 |
| Bank of New York Mellon | $1.4 | $0.9 | 12 | 22 | 74 |
| Citigroup | ($27.7) | $5.3 | 124 | 176 | 738 |
| Goldman Sachs | $2.3 | $4.8 | 212 | 391 | 953 |
| J P Morgan Chase | $5.6 | $8.7 | >200 | | 1626 |
| Merrill Lynch | ($27.6) | $3.6 | 149 | | 696 |
| Morgan Stanley | $1.7 | $4.5 | 101 | 189 | 428 |
| State Street Corp | $1.8 | $0.5 | 3 | 8 | 44 |
| Wells Fargo & Co. | ($42.9) | $1.0 | 7 | 22 | 62 |

*Source:* Cuomo (2009). Wells Fargo losses include losses from Wachovia (acquired in December 2008).

bonuses while losing $81.4 billion in earnings. Not surprising, the Cuomo report further fueled outrage over Wall Street bonuses on both Main Street and in Washington.

### 3.9.3 Treasury Issues Final Rules and Appoints a Pay Czar

The Dodd Amendments were signed into law with the understanding that the US Treasury would work out the implementation details. In June 2009, Treasury issued its rulings, along with the simultaneous creation of the Office of the Special Master of Executive Compensation. The Special Master (colloquially known as the Pay Czar) had wide-ranging authority over all TARP recipients, but was particularly responsible for all compensation paid to the top 25 executives in the seven firms deemed to have required special assistance from the US government: Bank of America, Citigroup, AIG, General Motors, Chrysler, and the financing arms of GM and Chrysler.[142]

Since taxpayers had become the major stakeholder in the seven special assistance firms, the government arguably had a legitimate interest in the firm's compensation policies. One could imagine, for example, embracing an objective of maximizing shareholder value while protecting taxpayers, or perhaps maximizing taxpayer return on investment. However, the US Treasury instructed the Special Master to make pay determinations using the "public interest standard", an ill-defined concept that allows too much discretion and destroys accountability for those exercising the discretion. For example, applying the public interest standard allows Congress to limit compensation they perceive as excessive, without evidence or accountability for the consequences.

[142] For the record, I (along with Lucian Bebchuk from Harvard) served as academic advisors to Kenneth Feinberg, the Special Master. However, that the fact advice was given does not imply that it was followed.

Similarly, invoking the public interest standard forced the Special Master to navigate between the conflicting demands of politicians (insisting on punishments) and tax-payer/shareholders (concerned with attracting, retaining, and motivating executives and employees). Ultimately, the Special Master catered to prevailing political and public sentiment, and severely penalized the executives in firms viewed as responsible for the meltdown by drastically reducing their cash compensation.

## 3.10  The Dodd–Frank Executive Compensation Reform Act (2010–2011)

In July 2010, President Obama signed into law the Dodd–Frank Wall Street Reform and Consumer Protection Act or Dodd–Frank Act, which was the culmination of the President and Congress's controversial and wide-ranging efforts to regulate the financial services industry. In spite of its enormous length—the bill itself spans 848 pages—the Act leaves most of the details to be promulgated by a variety of government entities. Indeed, attorneys at DavisPolk (2010) calculate that the Act requires regulators from at least nine agencies to create 243 new rules, conduct 67 studies, and issue 22 periodic reports.

### 3.10.1  Pay Restrictions for Financial Institutions

While the pay restrictions in the TARP legislation apply only to banks receiving government assistance, the Dodd–Frank Act goes much further by regulating pay for *all* financial institutions (public or private, TARP recipients and non-recipients) including broker–dealers, commercial banks, investment banks, credit unions, savings associations, domestic branches of foreign banks, and investment advisors. Specifically, Part (a) of Section 956 of the Dodd–Frank Act requires all financial institutions to identify and disclose (to their relevant regulator) any incentive-based compensation arrangements that could lead to material financial loss to the covered financial institution, or that provides an executive officer, employee, director, or principal shareholder of the covered financial institution with excessive compensation, fees, or benefits. In addition, Part (b) of Section 956 of the Dodd–Frank Act prohibits financial institutions from adopting any incentive plan that regulators determine encourages inappropriate risks by covered financial institutions, by (1) providing an executive officer, employee, director, or principal shareholder of the covered financial institution with excessive compensation, fees, or benefits; or (2) that could lead to material financial loss to the covered financial institution.

Since at least the early 1990s, there has always been a tension between shareholders (the firm's legal owners) concerned about CEO incentives, and uninvited guests (such as politicians and labor unions) concerned about high levels of pay. After the TARP bailouts in the financial crisis, the analogous tension was between taxpayers (who wanted to be protected from excessive risks while receiving an appropriate return on their investment) and politicians who were outraged about perceived excesses in banking bonuses. Section 956(b) of the Dodd–Frank Act deliberately conflates these tensions, by explicitly defining excessive compensation as an inappropriate risk. Moreover, Section 956(a)

of the Dodd–Frank Act requires banks to inform their regulators of compensation plans that provide excessive compensation, delegating to the regulators the Herculean task of defining what compensation is excessive (or, indeed, which risks are inappropriate).

The responsibility for implementing Section 956 of the Dodd–Frank Act fell jointly to seven agencies: the Securities and Exchange Commission (SEC), the Federal Reserve System, the Office of the Comptroller of the Currency, the Office of Thrift Supervision, the Federal Deposit Insurance Corporation, the National Credit Union Administration, and the Federal Housing Finance Agency. In March 2011, the seven agencies issued a joint proposal for public comment, modeled in part on Section 39 of the Federal Deposit Insurance Act. While the proposal stops short of explicitly limiting the level of executive compensation, it prohibits compensation that is unreasonable or disproportionate to the amount, nature, quality, and scope of services performed. In addition, the proposal calls for firms to identify individuals who have the ability the expose the firm to substantial risk, and demands that (for the larger institutions) such individuals have at least 50% of their bonuses deferred for at least three years; deferred amounts would be subject to forfeiture if subsequent performance deteriorates. Final rules were expected in late 2012.

### 3.10.2  Pay and Governance Reforms for all Publicly Traded Companies

While ostensibly focused on regulating firms in the financial services industry, the authors of the Dodd–Frank Act seized the opportunity to pass a sweeping reform of executive compensation and corporate governance imposed on all large publicly traded firms across all industries. The new rules include:

Say-on-Pay (Section 951). Shareholders will be asked to approve the company's executive compensation practices in a non-binding vote occurring at least every three years (with an additional vote the first year and every six years thereafter to determine whether the Say-on-Pay votes will occur every one, two, or three years). In addition, companies are required to disclose, and shareholders are asked to approve (again, in a non-binding vote), any golden parachute payments in connection with mergers, tender offers, or going-private transactions.

> In January 2011 – and effective for the 2011 proxy season – the SEC adopted rules concerning shareholder approval of executive compensation and "golden parachute" compensation arrangements. Shareholders of 98.5% of the 2532 companies reporting 2011 results by July 2011 approved the pay plans; over 70% of the companies received more than 90% favorable support.[143] Similarly, shareholders of 98.2% of the 1875 companies reporting 2012 results by June 2012 approved the pay plans; 72% of the companies received more than 90% favorable support.[144] Twenty six of the 30 companies receiving less than 50% positive votes in 2011

---

[143] Holzer, "A 'Yes' In Say On Pay," *Wall Street Journal* (2011b).

[144] "2012 Say on Pay Results" Semler Brossy Consulting Group, LLC. Accessed 8/6/2012 at www.semlerbrossy.com/sayonpay.

*passed in 2012, and year-over-year favorable votes increased by 14% for companies receiving between 50% and 70% favorable votes in 2011. Clawbacks (Section 954).*

Companies must implement and report policies for recouping payments to executive based on financial statements that are subsequently restated. The rule applies to any current or former executive officer (an expansion of Sarbanes–Oxley, where only the CEO and CFO were subject to clawbacks), and applies to any payments made in the three-year period preceding the restatement (Sarbanes–Oxley only applied for the twelve months following the filing of the inaccurate statement).

> *The SEC had neither adopted nor proposed rules regarding the recovery of executive compensation by August 2012. However, Equilar reports that 86% of the Fortune 100 companies issuing proxy statements in 2012 had publicly disclosed clawback arrangements; in half of the companies the clawback triggers were related to financial restatements and ethical misconduct.[145] Additional Disclosures (Sections 953, 955, 972).*

Companies must report the ratio of CEO compensation to the median pay for all other company employees. Companies must analyze and report the relation between realized compensation and the firms financial performance, including stock–price performance. In addition, companies must disclose its policies regarding hedging by employees to protect against reductions in company stock prices. Finally, the Dodd–Frank Act requires companies to disclose their policies and practices on why the company chooses either to separate the Chairman and CEO positions, or combine both roles.

> *The SEC had neither adopted nor proposed rules regarding the disclosure of pay ratios, pay-for-performance, hedging and CEO/Chair combinations by August 2012. Compensation Committee Independence (Section 952).*

Following Sarbanes–Oxley (2002) requirements for Audit Committees, publicly traded companies are required to have compensation committees comprised solely of outside independent directors (where independence takes into account any financial ties the outside directors might have with the firm. In addition, companies must assess the independence of compensation consultants, attorneys, accountants, and other advisors to the compensation committees.

> *In June 2012, the SEC adopted final rules directing exchanges to establish listing standards guaranteeing that members of the compensation committee (or directors who oversee executive compensation matters in the absence of a committee) to be independent. While leaving the precise definition of "independence" to the exchanges, the final rule required exchanges to consider the director's source of compensation (including consulting or advisory fees) paid by the issuer, and whether the director is affiliated with the issuer, a subsidiary of the issuer, or an affiliate of a subsidiary of the issuer.*
>
> *In addition, the new SEC rules require firms to ensure that compensation committees have authority and funding to retain compensation consultants. While neither the Act nor the June*

---

[145] *2012 Fortune 100 Clawback Report*, Equilar, Inc. August 2, 2012.

*2012 Final Rule issued by the SEC required compensation advisors to be independent, the SEC imposed a list of independence criteria that boards must consider in retaining a consultant. Finally, proxy statements issued in connection with annual shareholder meetings in 2013 and after must disclose whether the work of the consultant has raised any conflict if interest and, if so, the nature of the conflict and how the conflict is being addressed. Proxy Access (Section 971).*

The Dodd–Frank Act authorized the SEC to issue rules allowing certain shareholders to nominate their own director candidates in the company's annual proxy statements.

*The SEC issued its rules on Proxy Access in August 2010, but delayed implementation after lawsuits by the Business Roundtable and the US Chamber of Commerce claimed that the rules would distract management and advance special-interest agendas. In July 2011, the US Circuit Court of Appeals (Washington, DC) ruled in favor of the business groups and rejected the SEC's rule. As of August 2012, the SEC had not announced whether it would attempt to rewrite the rule in a way that would be acceptable to the Court.*

It is too early to assess the ultimate effect of Dodd–Frank on executive compensation, since many of the rules have just been implemented or are still being written. However, based on experiences with similar rules, I can speculate on the ultimate impact.

*Say on Pay.* In mandating Say-on-Pay, the Dodd–Frank Act follows similar rulings for non-binding shareholder votes enacted in the United Kingdom in 2002 and later in Australia, Denmark, France, Portugal, Spain, and Sweden; the Netherlands and Norway went a step further by allowing *binding* shareholder votes. Say-on-Pay had long been a favorite objective of Democrats in Congress, and the Say-on-Pay Bill passed the House in April 2007 by a 2:1 margin. While the companion bill introduced in the Senate by then-Senator Obama was shelved prior to a vote, Say-on-Pay was widely expected to become law following the 2008 presidential election, especially after Say-on-Pay was mandated for TARP recipients as part of the Dodd Amendments.

In spite of the support, however, there is modest evidence that Say-on-Pay results in important changes to compensation practices. In the United Kingdom (where we have the most data), there is some evidence that negative Say-on-Pay votes have led to some reductions in salary continuation periods in severance agreements and some changes in performance-based vesting conditions in equity plans, but no evidence that the votes have affected compensation levels (Ferri and Maber, 2010). In the United States, where shareholders voted on the compensation for TARP executives for the first time in early 2009, the plans were passed at all firms, with an average of 88.6% of the votes cast in favor of management. Among the TARP recipients garnering the strongest support were the Wall Street firms whose compensation systems allegedly fostered the financial crisis, including Goldman Sachs (98%), AIG (98%), JPMorgan (97%), Morgan Stanley (94%), Citigroup (84%), and Bank of America (71%).[146]

---

[146] Tse, "Shareholders Say Yes To Executive Pay Plans; Review Tracks Advisory Votes at TARP Firms," *Washington Post* (2009). It is worth noting that shareholders voting in early 2009 were largely voting on 2008 compensation, before the implementation of the Dodd Amendments or the appointment of the Special Master.

As emphasized in this chapter, regulation inevitably produces unintended consequences. The most obvious (and most negative) unintended consequence associated with Say-on-Pay reflects the increasing influence of proxy-advisory firms (primarily Institutional Shareholder Service (ISS)). To fulfill their required fiduciary duties to vote proxies, institutional investors routinely rely on ISS and other proxy-advisory firms for recommendations on how to vote on Say-on-Pay and other proxy matters. In turn, the proxy-advisory firms rely on a limited (and controversial) set of quantitative criteria to determine whether to offer positive or negative voting recommendations.[147] In a broad sample of Russell 3000 firms, Larcker, McCall, and Ormazabal (2012) show: (1) the recommendations of the proxy-advisory firms do, indeed, affect voting outcomes; (2) anticipating this result, firms change their compensation policies to avoid negative recommendations; and (3) the market reaction to these changes is statistically negative.

Firms inherently face different competitive and incentive challenges, and there is neither a "one-size-fits-all" solution to these challenges, nor a limited set of quantitative criteria that can substitute for a careful and holistic assessment of compensation plans that takes into account company-specific situations and objectives. Ultimately, the benefits of adhering to the ISS criteria must be weighed against the cost associated with reduced innovation and flexibility in the provision of compensation and incentives.

*Compensation Committee and Advisor Independence.* The Dodd–Frank provisions on the independence of the compensation committee will have little practical effect for large companies, since the listing requirements of the NYSE and NASDAQ have required independent compensation committees since 2003, and the IRS has required independent compensation committees (for Section 162(m) purposes) since 1994. The provision related to the independence of compensation consultants, in combination with SEC disclosure rules introduced in December 2009, will encourage more committees to retain their own independent consultant in addition to the consultants engaged by management.[148]

*Clawback Provisions.* The Sarbanes-Oxley experience shows that companies rarely try to recover erroneously awarded compensation from their CEO and CFO, often citing potential litigation costs and the feasibility of recouping money that has already been paid and taxed. The Dodd–Frank provision makes it more difficult for boards to shirk their responsibility to recovery erroneously awarded pay, and indeed likely subjects boards to shareholder litigation if they fail to even try.

---

[147] For critiques of the ISS methodology, see Wachtell-Lipton's "Say on Pay 2012," available at: http://blogs.law.harvard.edu/corpgov/2012/07/14/say-on-pay-2012/ (accessed 8/6/2012).

[148] The 2009 SEC disclosure rules require companies to disclose the fees paid to executive compensation consultants for any work beyond executive compensation (e.g. actuarial work, benefits administration, employee pay, etc.), but offers a safe harbor (i.e. no disclosure requirement) when the committee retains their own independent consultant. Interestingly, Murphy and Sandino (2010) find that levels of CEO pay are significantly higher in firms with consultants working exclusively for the compensation committee.

*Ratio of CEO-to-Worker Pay.* The most mischievous and controversial compensation provision in Dodd–Frank is the required disclosure of the ratio of CEO pay to the median pay of all employees. The calculation costs alone can be immense for large multinational or multi-segment corporations where payroll is decentralized: to compute the median the company needs an often non-existent single compensation database with all employees worldwide. More importantly, however, is what shareholders are supposed to do with this new information, or how they should determine whether a ratio is too high or too low. Ultimately, this provision reflects a belief in Congress that CEO pay is excessive and its sole purpose is the hope that disclosing the ratio will shame boards into lowering CEO pay.

*Proxy Access.* Finally, potentially most important is the Proxy Access rule allowing shareholders to include their director nominees on the proxy alongside with the board's nominees. In issuing its rule in August 2010, the SEC limited access to shareholders who have held at least 3% of the company's stock for at least three years. One view is that Proxy Access will provide shareholders with a critical mechanism to replace poor directors with better ones. A more-cynical view—expressed by the *Wall Street Journal* and others—is that 3% was chosen as the sweet spot for labor unions and other politically motivated organizations who will use their leverage over the proxy statement to force companies to support political causes rather than increasing shareholder value.[149] In its July 2011 ruling rejecting the SEC's rule, the US Circuit Court of Appeals (Washington, DC) issued a sharp rebuke to the SEC, saying that the SEC failed in analyzing the cost the rule imposes on companies and in supporting its claim that the rule would improve shareholder value and board performance.[150]

## 4. INTERNATIONAL COMPARISONS: ARE US CEOS STILL PAID MORE?

### 4.1 The US Pay Premium: What We Thought We Knew[151]

Among the best-known "stylized facts" about executive compensation is that CEOs in the United States are paid significantly more than similarly situated CEOs in foreign corporations (e.g. Abowd and Bognanno, 1995; Abowd and Kaplan, 1999; Murphy, 1999). However—although widely accepted by academics, regulators, and the media—this stylized fact has *not* generally been based on consistent and comprehensive pay data across a large number of countries with controls for cross-country differences in firm characteristics. In particular, while the United States has required detailed disclosures on executive compensation since the 1930s, the majority of other countries have historically required reporting (at most) the aggregate cash compensation for the top-management team, with no individual data and little information on the prevalence of equity or option grants.

---

[149] "Alinsky Wins at the SEC," *Wall Street Journal* (2010).
[150] Holzer, "Corporate News: Court Deals Blow to SEC, Activists," *Wall Street Journal* (2011a).
[151] This section draws heavily from Fernandes et al. (2013), Conyon et al. (2013).

In fact, prior to 2000, only Canada (which mandated pay disclosures in 1993) and the United Kingdom (based on disclosure recommendations issued in 1995) required US-style full disclosure of CEO compensation (including details on equity grants). Based on data from 1993 to 1995, Zhou (2000) shows that US CEOs earned more than double their Canadian counterparts. Conyon and Murphy (2000) show that US CEOs earned almost 200% more than British CEOs in 1997, after controlling for industry, firm size, and a variety of firm and individual characteristics. Conyon, Core, and Guay (2011) show that the US versus UK Pay Premium had fallen to 40% by 2003 and plausibly disappears after adjusting for the risk associated equity-based compensation.

Other multi-country pay comparisons have typically relied on aggregate or average executive pay across groups of executives, usually excluding equity-based pay).[152] For example, Conyon and Schwalbach (2000)'s comparison of UK and German compensation from 1968 to 1994 focused on only cash compensation for the United Kingdom (because the study predated the UK recommendations on disclosing stock options) and average cash compensation for Germany (because German rules required only disclosing the total cash paid across the group of top managers). Similarly, Muslu (2008)'s study of the largest 158 European companies from 1999 to 2004 (based on hand-collected annual reports) presents a mixture of individual and aggregated compensation data. Bryan, Nash, and Patel (2006) relied on SEC Form 20-F filings from 1994 to 2004 for foreign companies cross listing in the United States; however, cross-listed companies are only required to disclose compensation for individual executives if such disclosure is required in the home country, and as a result most of their analysis was based on average compensation for the management group.

Beyond the comparisons with Canada and the United Kingdom, and the handful of studies based on aggregate cash compensation data, much of what we know (or thought we knew) about international differences in CEO pay has been based on Towers Perrin's biennial *Worldwide Total Remuneration* reports, utilized (for example) by Abowd and Bognanno (1995), Abowd and Kaplan (1999), Murphy (1999), and Thomas (2008) (not coincidentally, the same cites as in the first paragraph). These international comparisons—which have typically suggested that US CEOs are paid more than twice the "going rate" for CEOs in other countries—are not based on "data" per se, but rather depict the consulting company's estimates of "typical" or "competitive" pay for a representative CEO in an industrial company with an assumed amount in annual revenues, based on questionnaires sent to consultants in each country. While crudely controlling for industry and firm size (by design), it is impossible using these surveys to control for other factors that might explain the US "pay premium", such as ownership and board structure or individual CEO characteristics.

---

[152] Single-country studies based on aggregate pay include Kaplan (1994) (Japan), Kato and Long (2005) (China), Fernandes (2008) (Portugal), and Kato, Kim, and Lee (2006) (Korea).

The disclosure situation has improved markedly over the past decade. Regulations mandating disclosure of executive pay were introduced in Ireland and South Africa in 2000 and in Australia in 2004. In May 2003, the European Union (EU) Commission issued an "Action Plan" recommending that all listed companies in the EU report details on individual compensation packages, and that EU member countries pass rules requiring such disclosure. By 2006, six EU members (in addition to the United Kingdom and Ireland) had mandated disclosure: Belgium, France, Germany, Italy, Netherlands, and Sweden. In addition, although not in the EU, Norway also adopted EU-style disclosure rules, and Switzerland demanded similar disclosure for the "highest-paid" executive.

## 4.2  New International Evidence

In my joint work with Nuno Fernandes, Miguel Ferreira and Pedro Matos (Fernandes et al., 2012)—based on recently available data from 14 countries with mandatory pay disclosures—we show that the stylized fact that US CEOs earn substantially more than foreign CEOs is wrong, or at least outdated. In particular, we show that the "US Pay Premium" became statistically insignificant by 2007 and largely reflects a risk premium for stock-option compensation (which remains more prevalent in the United States than in other countries).

In reaching our conclusion that the US Pay Premium has become modest (or insignificant), we control not only for the "usual" firm-specific characteristics (e.g. industry, firm size, volatility, and performance) but also for governance characteristics that systematically differ across countries: ownership and board structure. Compared to non-US firms, US firms tend to have higher institutional ownership and more independent boards, factors associated with both higher pay and increased use of equity-based compensation. In addition, shareholdings in US firms tend to be less dominated by "insiders" (such as large-block family shareholders), factors associated with lower pay and reduced use of equity-based compensation.

Figure 23 traces the evolution of the US pay premium from 2003 to 2008 (based on results in Table 8 of Fernandes et al., 2013). The premium is defined as $e^{\beta_1} - 1$ in the following regression, estimated annually for a pooled sample of US and non-US CEOs:

$$Ln(CEOPay_i) = \alpha + \beta_1(USDummy) + \beta_2(FirmCharacteristics_i) + \varepsilon_i$$

The sample consists of between 1426 and 1532 US firms and between 781 and 1480 non-US firms per year. US data are extracted from ExecuComp, while non-US data are based primarily on BoardEx and supplemented with hand-collected data from filings.

The "Firm Characteristics" in the left-hand panel of Figure 23 include only controls for company size (Ln(Revenues)) and industry (fixed effects for 12 Fama–French industries). As shown in the figure, the implied US Pay Premium fell significantly from over 100% in 2003–2005 to less than 80% in 2006–2008. The right-hand panel includes additional controls for leverage, Tobin's Q, stock volatility, stock returns, ownership structure

**Figure 23** The evolving (and disappearing) US Pay Premium. *Note:* The figure shows the US Pay Premium implied by regression Ln(CEO Pay) on a US dummy variable plus controls for industry and company revenues (left-hand panel) and also other firm characteristics, ownership structure, and board structure (right-hand panel) for each year from 2003 to 2008. ***, **, * indicates that the coefficient on the US premium on each underlying regression depicted above is significant at the 1, 5, and 10% levels, respectively. *Source:* Fernandes et al. (2013), Table 8.

(the fraction of shares held by insiders and institutions) and board structure (board size, independence, the average number of board positions held by each board member, and a dummy variable indicating that the CEO also holds the title of Chairman). As shown in the figure—after including these additional controls—the implied US Pay Premium declined from nearly 60% in 2003 to only 26% in 2006 and 2% in 2007.

Figure 24 shows the international distribution of predicted 2006 CEO pay for a hypothetical firm with $1 billion sales, based on the specification used for Figure 23 with the "US dummy" replaced by a set of 14 country dummies. Panel A, in the spirit of the Towers Perrin estimates, controls only for firm size and industry, while Panel B controls for industry, firm characteristics, ownership, and board characteristics. The pay composition percentages are defined as the average composition across all CEOs for each country. The figure shows that US CEOs earn substantially more than non–US CEOs controlling only for size and industry. However, after controlling for firm, ownership, and board characteristics, we find effective parity in CEO pay levels among Anglo–Saxon nations (United States, United Kingdom, Ireland, Australia, and Canada) and also Italy.

As an extension to the results in Figures 23 and 24, we also compare international differences in risk–adjusted pay, using methodologies similar to that used above in Section 2.1.2 and Figure 7.[153] Consistent with the conclusions of Conyon et al.

---

[153] Due to limitations with BoardEx data on CEO wealth for non–US CEOs, Fernandes et al. (2013) make simplifying assumptions beyond those in Section 2.1.2.

**Figure 24** 2006 CEO pay after controlling for firm characteristics, ownership, and board structure Panel A. Controlling only for sales and industry Panel B. Controlling for sales, industry, and firm, ownership, and board characteristics Note: The figure compares 2006 CEO pay in each country controlling for firm size (sales) and industry in Panel A, and controlling for size, industry, and firm, ownership, and board characteristics in Panel B. We regress the logarithm of total compensation on the logarithm of sales and 12 industry and 14 country dummies. For each country, we estimate the pay for a CEO running a hypothetical firm with $1 billion in sales using the estimated coefficient for pay-size sensitivity and controlling for the "average" industry. The non- US average is weighted by the number of firms in each country. The pay composition percentages are defined as the average composition across all CEOs for each country. *Source:*Fernandes, et al. (2013), Figure 1.

(2013) (who use a different methodology and consider only US–UK comparisons), we find that the risk-adjusted US pay premium for 2006 is statistically insignificant after controlling for governance (but remains significant before such controls), and that risk-adjusted pay in the US is significantly less than CEO pay in the United Kingdom and Australia, and insignificantly different from CEO pay in Canada, Italy and Ireland.

In addition, we show that both the level and structure of 2006 pay for US CEOs is insignificantly different from that of non-US CEOs of "internationalized" firms, which we define as firms above the 75th percentile ranked by foreign institutional ownership, foreign sales (as a fraction of total sales), or board international diversity (defined as the number of different nationalities represented on the board of directors divided by the total board size). We also find insignificant differences between US CEOs and non-US CEOs in firms included in the 1500-firm Morgan Stanley Capital International All Country World Index (routinely used as a benchmark for global equity mutual funds and used here as a proxy for foreign investor demand).

Finally, we find no significant differences in the level or structure of pay when US CEOs are compared to non-US CEOs of "Americanized" firms, which we define as firms cross-listed on US exchanges (as a proxy for demand by US investors) or above the 75th percentile ranked by US institutional ownership, total acquisitions of US companies between 1996–2005 (as a proxy for exposure to US product and labor markets), and the fraction of directors who also sit on boards of companies headquartered in the United States (as a proxy for exposure to US pay practices).

Overall, our evidence is inconsistent with the view that US CEO pay is "excessive" when compared to that of their foreign counterparts, but rather reflects tighter links between CEO pay and shareholder performance in US firms. First, we show that the US pay premium is modest after controlling for firm, ownership, board, and CEO characteristics. Second, we demonstrate that it is misleading to examine cross-sectional or cross-country differences in the *level* of pay in isolation, without also examining differences in the *structure* of pay, namely the use of equity-based compensation. In fact, the firm, ownership, and board characteristics associated with higher pay are those associated with a larger fraction of equity-based pay. Third, we find that CEO pay levels and the use of equity-based compensation are positively related to variables routinely used as proxies for better monitoring and better governance, namely institutional ownership and board independence. Fourth, our findings suggest that the observed US CEO pay premium reflects compensating differentials for the equity-based pay increasingly demanded by internationally diverse boards and shareholders. We find evidence that foreign and US institutional shareholders are linked to a greater use of equity-based pay and higher pay levels in non-US firms in which they invest. Finally, the convergence of US and non-US CEO pay levels since 2003 seems to be explained by the convergence of ownership structures and globalization of capital markets.

## 4.3 Why Do US CEOs Receive More Options?

Our finding that the US pay premium largely disappears after controlling for the relative riskiness of US pay packages potentially "explains" the pay differences but naturally leads to another question: Why do US executives receive more equity-based compensation than their foreign counterparts?

While equity-based compensation has been a staple of US compensation contracts for more than a half-century, the use of equity-based pay outside the United States is a relatively recent phenomenon. Panel A of Table 3 shows how the importance of equity-based pay has changed over time in the United States and in nine European countries using Towers Perrin's *Worldwide Total Remuneration* (WWTR) surveys for the selected

**Table 3** Stock-based pay (as a percentage of total pay) in Europe and the United States
**Panel A: Towers Perrin consultant surveys 1984–2003**

|                | 1984  | 1988  | 1992  | 1996  | 1999  | 2001  | 2003  |
|----------------|-------|-------|-------|-------|-------|-------|-------|
| Belgium        | 0.0%  | 0.0%  | 0.0%  | 0.0%  | 3.2%  | 11.6% | 11.2% |
| France         | 12.3% | 13.3% | 15.6% | 14.6% | 14.3% | 15.1% | 16.0% |
| Germany        | 0.0%  | 0.0%  | 0.0%  | 0.0%  | 9.7%  | 13.5% | 18.0% |
| Italy          | 0.0%  | 0.0%  | 0.5%  | 4.0%  | 9.1%  | 17.2% | 15.1% |
| Netherlands    | 0.0%  | 0.0%  | 0.0%  | 0.0%  | 14.6% | 16.7% | 15.8% |
| Spain          | 0.0%  | 0.0%  | 0.0%  | 0.0%  | 16.0% | 17.9% | 19.2% |
| Sweden         | 0.0%  | 0.0%  | 0.0%  | 0.0%  | 6.8%  | 11.0% | 10.7% |
| Switzerland    | 1.9%  | 1.9%  | 3.4%  | 3.6%  | 1.8%  | 0.0%  | 19.2% |
| United Kingdom | 14.5% | 14.6% | 15.7% | 15.0% | 16.6% | 19.1% | 20.8% |

**Panel B: BoardEx (non-US firms) and ExecuComp (US firms)**

|                | United States | 16.9% | 28.3% | 32.3% | 28.7% | 25.5% | 44.8% | 48.3% |
|----------------|---------------|-------|-------|-------|-------|-------|-------|-------|

|                | 2003  | 2004  | 2005  | 2006  | 2007  | 2008  |
|----------------|-------|-------|-------|-------|-------|-------|
| Belgium        | na    | 16.7% | 8.6%  | 9.5%  | 7.7%  | 11.5% |
| France         | 17.6% | 15.9% | 16.0% | 17.2% | 17.8% | 13.9% |
| Germany        | 12.5% | 8.7%  | 9.4%  | 9.3%  | 9.4%  | 9.0%  |
| Italy          | 11.5% | 10.6% | 15.7% | 13.1% | 5.7%  | 8.6%  |
| Netherlands    | 19.3% | 16.3% | 20.1% | 21.7% | 18.2% | 15.8% |
| Spain          | 0.0%  | 1.2%  | 0.0%  | 0.8%  | 5.3%  | 2.9%  |
| Sweden         | 3.7%  | 1.3%  | 1.5%  | 1.8%  | 1.5%  | 1.3%  |
| Switzerland    | 30.2% | 21.5% | 20.1% | 26.6% | 17.1% | 12.0% |
| United Kingdom | 27.7% | 27.7% | 29.7% | 31.0% | 34.1% | 30.6% |
| United States  | 40.7% | 42.0% | 41.4% | 39.5% | 43.0% | 47.1% |

*Note:* Data in Panel A are from Towers Perrin's Worldwide Total Remuneration reports (various issues), including 1984–1992 data reported by Abowd and Kaplan (1999). Data reflects Towers Perrin's estimate of competitive CEO pay for industrial companies with approximately US $300 million in annual revenues. Stock-based pay includes the grant-date expected value of option grants and annualized targets from performance share plans. Data in Panel B are from BoardEx and ExecuComp. The percentages in Panel B are constructed by first computing the average ratio of equity-based pay to total compensation for each CEO, and then averaging across CEOs.

years 1984, 1988, 1992, 1996, 1999, 2001, and 2003. The data for the years 1992–1996 are based on the Abowd and Kaplan (1999) analysis of the WWTR surveys. As shown in Panel A, only France and the UK made extensive use of stock or options in the 1980s, and equity-based pay did not become common across Europe until the end of the 1990s. By 2003, Towers Perrin reports that equity-based pay accounts for between 10% and 20% of competitive pay for European CEOs, and for about half the pay of American CEOs.

As discussed earlier, the data in Panel A of Table 3 are not CEO pay "data" per se, but rather consulting company's estimates of "typical" or "competitive" pay for a representative CEO in an industrial company, based on questionnaires sent to consultants in each country. In Panel B of Table 3, I provide my own estimates of equity-based pay for 2003–2008 based on actual grant-date values extracted from BoardEx (for Europe) and ExecuComp (for the United States). The actual averages for 2003 in Panel B are generally consistent with the consultant surveys in Panel A for the same year, increasing my confidence in both data sources. As shown in Panel B, the use of equity-based compensation has generally declined in continental Europe between 2003 and 2008, and has remained relatively constant in the United Kingdom at just under a third of total compensation. In contrast, the use of equity-based pay has increased in the United States.

Traditional agency theory suggests a finite number of factors that might explain a greater use of incentive-based pay among US executives. First, US CEOs may be less risk averse or have steeper marginal costs of effort than their non-US counterparts, but to our knowledge there is no theory or empirical work suggesting such international differences in risk-aversion coefficients. Second, performance of non-US firms might be measured with substantially more noise than for US firms, leading to lower pay-performance sensitivities and lower expected levels of pay. However, we find no evidence that cash flows or shareholder returns are systematically more variable in our sample of non- US firms than in US firms. Extensions of the traditional model to incorporate differences in both ability and in the marginal productivity of CEO effort might help reconcile the data, but only given the additional assumptions that executives are more able and more productive in the United States. Overall, there are no compelling agency-theoretic explanations for the relative reliance on equity-based compensation in the United States.[154]

In unreported analysis, we attempt to explain international differences in the use of equity-based compensation by a variety of country-level variables routinely used in international studies of corporate governance to measure differences in the economic, law, and institutional environment of each country.[155] We find that CEO equity-based pay (and

[154] Yermack (1995) shows that agency-theoretic variables have little explanatory value in predicting the use of equity-based compensation in a cross-section of US publicly traded firms.

[155] The limited number of countries in our sample (14) limits the statistical degrees of freedom for reliably identifying country-level determinants of pay practices.

total pay) is more prevalent in common-law countries (La Porta et al., 1998) which in turn is largely defined by the United Kingdom and its former colonies, including (in our sample) Australia, Canada, Ireland, South Africa, and the United States, and countries with stronger investor protections and private control of self-dealing (Djankov et al., 2008). We also consider different aspects of a country's regulatory environment. We find a positive association between CEO equity-based pay and the levels of compensation disclosure and director liability (La Porta, Lopez-De-Silanes, and Shleifer, 2006); note that the United States scores high in both indices. We find that equity-based pay is lower in countries with friendlier collective labor laws and countries where labor unions are more powerful (Botero et al., 2004), such as in Continental European countries (e.g. France and Germany). In contrast, differences in CEO pay are not explained by GDP per capita levels.

Ultimately, the cross-country differences in the prevalence of equity-based compensation may be driven by idiosyncratic events that in some cases encouraged, and in others discouraged, the use of stock options and restricted stock. For example, as documented in Section 3, America's reliance on stock options as the primary form of long-term compensation began in the 1950s as a result of tax policies designed to promote options, and declined in the late 1960s when the government reduced tax benefits. The early 1990s created a "perfect storm" for an explosion of option grants for not only executives but also lower-level managers and employees. The explosion in option grants continued unabated until the burst of the Internet bubble in 2000, followed by a series of accounting scandals that re-focused attention on the accounting treatment of options. Eventually, FASB mandated expensing, and companies moved away from options toward restricted stock.

Conyon et al. (2013) provide an analogous description of the evolution of equity-based pay in Europe. For example, the widespread adoption of stock-option plans in Europe initially emerged as governments provided tax incentives to encourage their use in the United Kingdom (in 1982), France (1984), and Italy (1998). Controversies in the United Kingdom in the 1990s involving perceived option excesses at recently privatized utilities led to a shift from options to restricted stock; concerns over excessive executive pay led France to revoke its tax subsidies on options in 1995, and Italy to revoke its tax subsidies in 2006. In Germany, option plans were not even legalized until 1996, and were still challenged in a series of high-profile lawsuits brought by a maverick college professor. In 1999, the Spanish government increased taxes on stock options after it was revealed that the CEO of the recently privatized telephone company was about to make a fortune exercising options.

In each country, ebbs and flows in option grants followed government intervention, usually reflecting tax or accounting policies and often reactions to isolated events or situations. Since the triggering events vary across countries, the nature of the government intervention—and the subsequent use of stock options—has also varied. The "perfect storm" that triggered the US option explosion (i.e. the "six factors" explored in Section 3.7 above) has not been repeated elsewhere in the world, and therefore the use of options (and equity-based pay in general) continues to be much higher in the United States.

## 5.  TOWARDS A GENERAL THEORY OF EXECUTIVE COMPENSATION

The academic literature focused on explaining cross-sectional differences and time-series trends in executive compensation is roughly divided into two camps: the "efficient contracting" camp and the "managerial power" camp. The efficient contracting camp maintains that the observed level and composition of compensation reflects a competitive equilibrium in the market for managerial talent, and that incentives are structured to optimize firm value. The managerial-power camp maintains that both the level and composition of pay are determined not by competitive market forces but rather by powerful CEOs, often working through or influencing captive board members. Most papers in the literature have adopted one approach or the other (often implicitly), and an increasing number of papers have treated the two approaches as competing hypotheses, attempting to distinguish between them empirically.

Ultimately, viewing efficient contracting and managerial power as competing hypotheses to "explain" executive compensation has not been productive. First, the hypotheses are not mutually exclusive; indeed, the same institutions that have evolved to mitigate conflicts of interest between managers and shareholders (i.e. efficient contracting) have simultaneously allowed executives to extract rents (i.e. managerial power). For example, the first "line of defense" against agency problems are the outside members of the board of directors, elected by shareholders and responsible for monitoring, hiring, firing, and setting top-executive compensation. However, these outside board members—who pay executives with shareholder money and not their own—are in no sense perfect agents for the shareholders who elected them. Instead of viewing efficient contracting and managerial power as competing hypotheses, it is more productive to acknowledge that outside-dominated boards mitigate agency problems between managers and shareholders but create agency problems between shareholders and directors. Rigidly adopting either extreme hypothesis—that director incentives are fully aligned with shareholder preferences or with those of incumbent CEOs—will inevitably result in less interesting and less realistic conclusions.

More importantly, viewing executive compensation as a "horse race" between efficient contracting and managerial power ignores other forces that may be even more important in explaining trends in pay. A central theme of this study is that government intervention into executive compensation—largely ignored by researchers—has been both a response to and a major driver of time trends in CEO pay. The reason political influence on CEO pay adds an important new dimension to the agency problem is because the interests of the government differ significantly from those of shareholders, directors, and executives. In particular, Congressional (and, more generally, populist) outrage over executive pay is almost always triggered by perceived excesses in the *level* of compensation without regard to incentives and company performance, and the regulatory responses have also fixated on pay levels (albeit with little effect). In contrast, while

shareholders have a legitimate concern over pay levels, their primary concern is whether executives have incentives to take actions that increase firm value, while avoiding value-destroying actions. Self-interested CEOs naturally prefer higher pay to lower pay. Directors, who are elected by shareholders but often selected by CEOs, appear to prefer better-aligned incentives but are not particularly interested in restraining pay levels.

## 5.1  Agency Problems: Solutions and Sources

The early 1900s witnessed the emergence of large publicly traded corporations with complex management structures that competed with and often displaced owner-managed and family-founded enterprises. Accompanying the rise in the "American Corporation" was the emergence of "professional executives"—non-owners hired to manage the firm's assets on behalf of passive and dispersed owner-shareholders (Wells, 2010). As noted by Smith (1776) in the context of 18th-century British "joint-stock" companies:

> *"Being managers rather of other people's money than of their own, it cannot well be expected, that they should watch over it with the same anxious vigilance with which the partners in a private copartnery frequently watch over their own . . . Negligence and profusion, therefore, must always prevail, more or less, in the management of the affairs of such a company."*

The conflicts identified by Smith (1776) arising between the owners of large publicly traded corporations and their hired executives is the quintessential "agency problem" explored by Berle and Means (1932) and Jensen and Meckling (1976). There are at least three versions of this agency problem:

• *The Agency Cost of Equity,* reflecting the fact that executives who own less than 100% of the shares of an all-equity firm will not make the same decisions (or "watch over it with the same anxious vigilance") they would if they owned 100% of the shares. Executives (usually assumed to be risk averse) want to be paid more and to take actions that increase their own utility, while shareholders (usually assumed to be risk neutral, or close to it) are primarily concerned with providing executives with incentives to take actions that increase the value of their shares.

• A variant of the Agency Cost of Equity is the "Agency Cost of Free Cash Flow" proposed by Jensen (1986a), reflecting the conflict of interest between executives and financial claimants on the disposition of cash flows in excess of those required to fund all positive net-present-value projects. While value is maximized by returning free cash flow to shareholders in the form of dividends or repurchases, empire-building executives prefer to retain and reinvest free cash flow unproductively in projects that destroy shareholder value. Debt financing mitigates free cash flow problems by pre-committing executives to pay out rather than retain future cash flows.

• *The Agency Cost of Debt,* reflecting the potential conflict of interest that exists between a company's shareholders and its debtholders: shareholders in a leveraged

firm prefer riskier investments than those that would maximize firm value, while debtholders prefer safer investments than those that would maximize firm value.[156] In addition, dividends and other payouts to shareholders may harm debtholders by jeopardizing the company's ability to service its debt. While the Agency Cost of Debt is clearly valid conceptually, there is little empirical evidence that leverage indeed leads to excessive risk taking, for several reasons. First, precisely because these conflicts are well understood, the potential problem is mitigated through debt covenants and constraints on how the proceeds from debt financing can be used. Moreover, since the problem is "priced" into the terms of the debt (with debtholders charging higher interest råtes in situations where executives have incentives to take higher risks), firms anticipating repeat trips to the bond market are directly punished for their risky behavior. The potential for conflicts are exacerbated, however, when the debtholders (or other fixed claimants, such as depositors) are protected against losses by the government. Such government guarantees can be explicit (such as FDIC insurance on deposits) or implicit (such as "Too Big To Fail" (TBTF) guarantees)). In these situations, the debtholders (or depositors) have little incentive to monitor management or enforce debt covenants, since the government is expected to cover losses.

While the labels on the various agency problems may be useful, they are all examples of the underlying agency conflict that arises when decision makers do not bear 100% of the wealth consequences of their decisions. As emphasized by Jensen (1993) there are four forces that mitigate agency problems between executives and the owners of large publicly traded corporations: (1) boards of directors; (2) capital markets; (3) the legal/political/regulatory system; and (4) product markets. However, while each of these forces can (and have) played a productive role in reducing agency conflicts, they also can (and have) created new problems, as follows.

### 5.1.1  Boards of Directors

The first line of defense against agency problems is the board of directors, elected by shareholders and responsible for monitoring, hiring, firing, and setting the compensation of the CEO and top-management team. For most of the prior century, boards were dominated by current executives and other corporate insiders. However, beginning with the shareholder movement in the 1980s (Section 3.6.2 above), firms have faced pressures for increased outsider representation on boards. By the end of the 1990s, the fraction of outside directors serving on the average board had increased to 80%, and the CEO was the sole insider in nearly half of all firms (Horstmeyer, 2011).

---

[156] As discussed in Section 2.2.2, it is not leverage *per se* that creates risk-taking incentives, but rather the limited liability feature of equity. The severity of the risk-taking incentives depends on the maximum downside risk compared to the dollar amount of equity, and not the value of equity compared to the overall value of the firm.

Conceptually, outside directors reduce agency problems by threatening errant executives with termination and by implementing incentive contracts that tie pay to value creation. The contracts that evolve from this setting will typically tie CEO pay to the creation of shareholder value, thus providing the theoretical justification for stock options, restricted stock, and other forms of equity-based compensation. Under the efficient contracting hypothesis, the contracts will be those that maximize shareholder value, while paying the CEO enough "expected" compensation or utility to get him to take the job, and recognizing that CEOs will respond predictably to the incentives provided by the contract.[157]

However, outside directors—who often own only a trivial fraction of their firm's common stock—are in no sense perfect agents for the shareholders who elected them. Board members are "reluctant to terminate or financially punish poor-performing CEOs because [board members] personally bear a disproportionately large share of the non-pecuniary costs [of such terminations], but receive essentially none of the pecuniary benefits" (Baker, Jensen, and Murphy, 1988, p. 614). Similarly, board members are willing to over-compensate adequately performing CEOs, since they are paying with shareholder money and not their own. As documented by Fracassi and Tate (2012), even "outside" board members often share important social ties with incumbent CEOs, especially in cases with powerful CEOs who presumably influence the director-nomination process. This agency problem between shareholders and their elected representatives forms the basis of the "managerial power hypothesis", in which powerful CEOs are able to influence both the level and composition of their own compensation packages. However, as discussed in Section 5.2.1 below, the agency problem is perhaps even more apparent in situations not involving powerful incumbents, such as directors overpaying CEOs hired from the outside.

### 5.1.2  Capital Markets

As discussed in Section 3.6.1, the executive compensation practices of the 1970s provided few incentives for executives to pursue value-increasing reductions in excess capacity and disgorgements of excess cash. However, pressures to improve performance, disgorge cash, and create wealth were ultimately introduced by the capital markets. The takeovers in the 1980s—often financed with newly available high-yield debt—provided credible competition for poorly performing incumbent managers. Wealth was created by both the post-merger activities of the acquiring firms (such as firing incompetent incumbent managers) and by responses to the takeover threat (such as excess spending cash to repurchase shares). Debt created value by providing commitments that the firm

---

[157] The optimal-contracting or principal-agent theory evolved contemporaneously to, but largely independent from, the agency-theory literature spawned by Jensen and Meckling (1976). Influential early theoretical work includes Ross (1973), Mirrlees (1976), Holmstrom (1979), Lazear and Rosen (1981), Holmstrom (1982), and Grossman and Hart (1983).

would pay its cash flows to debtholders, reducing the amounts available for executives to waste.

Capital markets—in particular, shareholder activists and large-block institutional stockholders—have mitigated agency problems by pressuring companies to strengthen links between CEO wealth and company stock-price performance. Fernandes et al. (2013), for example, show that the fraction of CEO pay delivered in the form of stock or options increases with institutional ownership. Hartzell and Starks (2003) show that CEO pay-performance sensitivities increase with the concentration of institutional ownership. In an international study, Aggarwal et al. (2011) find that the performance-related CEO turnover also increases with institutional ownership.

Capital markets have also, however, contributed to agency problems by providing executives with incentives to take actions to meet or beat analyst and market earnings expectations. As discussed in Section 2.4 and shown in Figure 12, executives have incentives to beat analysts forecasts by a small amount but not by *too much* because the abnormal stock-price response from beating the forecast by a lot is not much higher than the response for beating it by a little. Moreover, if an executive is going to miss the forecast, the executive may as well miss it by a *lot* since the negative abnormal stock-price response for a large miss is not much higher than for a small miss.

More generally, as argued by Jensen and Murphy (2012) and Martin (2011), capital-market pressures teach executives to focus on the "expectations market" (in which investors bet on expectations of future performance) rather than the "real market" (in which goods and services are produced and sold, and value is created or destroyed). Focusing on the expectation market is problematic because executives inherently have access to information about future prospects that are not publicly known and incorporated into stock prices. Executives with such a focus will be tempted to take actions that increase short-run stock prices at the expense of long-run value.

Temptations to manipulate the expectations market will clearly be higher for executives holding large quantities of stock and options that can be sold or exercised before markets adjust to the "real" information. As discussed in Section 2.4, there is substantial evidence that executive option and equity holdings are indeed higher in companies that restate their earnings or are accused of accounting fraud.[158] There is less evidence, however, that executives actually exercise and sell large fractions of their exercisable options or sell large fractions of their unrestricted stock holdings prior to restatements or indictments. The ominous hypothesis is that executives focused on the expectations market are not following a "pump and dump" strategy (which can be controlled by imposing longer holding requirements for shares), but rather that they are legitimately confused about the difference between increases in the short-run stock price and true value creation.

---

[158] See, for example, Efendi, Srivastava, and Swanson (2007), Burns and Kedia (2006), Bergstresser and Philippon (2006), Johnson, Ryan, and Tian (2009), and Erickson, Hanlon, and Maydew (2006).

### 5.1.3  The Political, Legal, and Regulatory System

Agency costs are mitigated by laws prohibiting embezzlement, corporate theft, and fraudulent conveyance, as well as securities rules, regulations, and listing requirements designed to protect shareholders and other financial claimants. For example, the Securities Act of 1933—which regulated new securities issues—sought to protect shareholders by mandating full disclosure of all information that a "reasonable shareholder" would require in order to make up his or her mind about the potential investment. The Securities Act of 1934—which regulated secondary trading of securities—introduced in Section 16(b) the "short swing" profit rule (discussed above in Section 3.5.4) requiring executives to return any profits realized from buying and selling (or selling and buying) shares of their company's stock within any period of less than six months. More sweeping (at least in its interpretation) was the anti-fraud provision Section 10(b) (and the corresponding SEC 10b-5 rule), which restricts insider trading, earnings manipulation, and price fixing. More recently, Regulation FD (August 2000) requires publicly traded companies to disclose material information to all investors at the same time (rather than favoring certain investors). While there are substantive arguments for allowing trading on material nonpublic information (since new information is more quickly introduced into the market), insider-trading rules are generally believed to benefit shareholders by reducing self-dealing by unscrupulous executives.

In addition to the general Securities Acts, the government has directly regulated the composition of the board of directors. Since 1994, companies have been required to have compensation committees consisting solely of independent directors in order for any pay to be exempt from the $1 million deductibility cap. In 1999, full independence of the auditing committee was required for all NYSE-listed firms; this requirement was extended to all firms in the 2002 Sarbanes-Oxley Act. In 2003, NYSE and NASDAQ listing requirements tightened the definition of independence and mandated that boards of listed firms have a majority of outside directors; the NYSE further required full independence for the compensation and nominating committees.

Critics hoping that independence requirements would reduce levels of executive pay have been disappointed. Both the level of pay and the use of equity-based compensation increase with the fraction of outsiders on the board; Fernandes et al. (2013) show that pay levels increase with board independence even after controlling for the risk associated with higher incentives. The evidence is therefore consistent with the hypothesis that directors—paying with shareholder money and not their own—prefer better-aligned incentives but are not particularly interested in restraining pay levels. The evidence is also consistent with directors not fully understanding (or believing) the opportunity cost of equity-based compensation (see Section 5.2.3 below).

Moreover, evidence that board independence "improves" pay is elusive. Bizjak and Anderson (2003) analyze the level and structure of compensation for CEOs who sit on their companies' compensation committees (a relatively common occurrence before the

early 1990s). Most critics of CEO pay (including Bebchuk–Fried and many shareholder activists) are horrified by the idea that the CEO could be a member of his own compensation committee, and would predict that such CEOs would inflate their own pay with few constraints.[159] And yet, Bizjak and Anderson (2003) find that the CEOs sitting on their own compensation committees earn substantially less (and not more) than other CEOs, have significant shareholdings and are typically company founders or their family members. These CEOs sit on their compensation committees not to inflate their own salaries, but rather to influence the level and structure of pay for their subordinates. Prohibiting such CEOs from sitting on (or chairing) their compensation committees harms shareholders, and illustrates a cost of the "one-size-fits-all" nature of corporate governance regulation.

In addition to general securities laws and independence requirements, this study has chronicled the history of government intervention into executive compensation. Over the past 80 years, Congress has imposed tax policies, accounting rules, disclosure requirements, direct legislation, and other rules designed explicitly to address perceived abuses in executive compensation. With few exceptions, the regulations have been either ineffective or counterproductive, typically increasing (rather than reducing) agency problems and pay levels, and leading to a host of unintended consequences. For example, the 1984 laws introduced to reduce golden parachute payments led to a proliferation of change-in-control arrangements, employment contracts, and tax gross-ups. Similarly, a variety of rules implemented in the early 1990s are largely responsible for fueling the subsequent option explosion, and the enhanced disclosure of perquisites in the 1970s is generally credited with fueling an escalation in the breadth of benefits offered to executives.

The emerging conclusion is that the myriad attempts to regulate CEO pay have been mostly unblemished by success. Part of the problem is that regulation—even when well intended—inherently focuses on relatively narrow aspects of compensation allowing plenty of scope for costly circumvention. An apt analogy is the Dutch boy using his fingers to plug holes in a dike, only to see new leaks emerge. The only certainty with pay regulation is that new leaks will emerge in unsuspected places, and that the consequences will be both unintended and costly.

Another part of the problem—as suggested above in the context of CEOs sitting on their firm's compensation committees—is that government regulation inevitably imposes a "one-size-fits-all" solution to a perceived problem. For example, as I emphasize in Murphy (2012), claims (unfounded or not) that the banking bonus culture created incentives to take excessive risks were relevant at most for a relatively small number

---

[159] While it was relatively common for CEOs to sit on their own compensation committees, I am unaware of any instances where the CEO was actually allowed to vote on his or her individual compensation package.

of large publicly traded Wall Street security brokers and dealers (along with some large commercial banks with significant investment banking operations). And yet, the Dodd–Frank provisions designed to reduce such incentives in the future were imposed on all public and private financial institutions, including broker-dealers, commercial banks, investment banks, credit unions, savings associations, domestic branches of foreign banks, and investment advisors.

A larger part of the problem is that the regulation is often mis-intended. The regulations are inherently political and driven by political agendas, and politicians seldom embrace "creating shareholder value" as their governing objective. While the pay controversies fueling calls for regulation have touched on legitimate issues concerning executive compensation, the most vocal critics of CEO pay (such as members of labor unions, disgruntled workers and politicians) have been uninvited guests to the table who have had no real stake in the companies being managed and no real interest in creating wealth for company shareholders. Indeed, a substantial force motivating such uninvited critics is one of the least attractive aspects of human beings: jealousy and envy. Although these aspects are seldom part of the explicit discussion and debate surrounding pay, they are important and impact how and why governments intervene into pay decisions.

### 5.1.4  Product Markets

While competition in the product market can theoretically either reduce or increase agency problems (see Hart, 1983; Scharfstein, 1988 respectively), companies that cannot compete in the product market cannot survive. The product market, therefore, provides inevitable discipline for value-destroying managers, but only after most of the value has been destroyed. Moreover, relying on product markets to discipline managers encourages managers to view "survival" rather than value-creation as their governing objective.

## 5.2  "Competing" Hypotheses to Explain the Increase in CEO Pay

The unparalleled rise in CEO pay from the mid-1980s through 2001—propelled primarily by increases in the grant-date value of option awards—generated a great deal of academic, popular, and political attention. As noted, most papers in the literature have offered either the "managerial power" or "efficient contracting" explanations for the increase; see Frydman and Jenter (2010) for a useful and thoughtful review. A third set of explanations—most closely associated with Murphy (2002)—maintains that options exploded in the 1990s because decisions over options were made based on the "perceived cost" of options rather than on their economic cost. This section summarizes and critiques all three approaches, focusing on salient features of CEO pay that can, and cannot be explained under the approach. In addition, I explore the government's role in pursuing social policy that favored stock options for both top-level executives and lower-level employees.

Before assessing how well the various theories explain the recent trends in CEO pay, it is useful to summarize what those trends are (that is, what the theories need to explain):

- Median expected pay for CEOs in the S&P 500 increased an average of 4.3% annually (after inflation) from 1983-1991, and by an average of 15.7% annually between 1991 and 2001.
- Most of the increase in pay between 1991 and 2001 reflects increases in the value of stock options granted.
- The "stock option explosion" was not limited to CEOs: 95% of the option grants went to lower-level executives and employees, and the trends in CEO options mirrored trends for options to lower levels.
- Median CEO pay has largely leveled-off since 2001. Over the same time period, firms have reduced their reliance on stock options and greatly increased their use of restricted stock and performance shares.

Therefore, any compelling theory of trends in CEO compensation must not only explain the increase in pay levels but must also address explicitly its most prominent feature: the escalation in stock options from the mid-1980s through 2001. Better still, the theory should be consistent with the explosion in broad-based option programs, the leveling of pay after 2001 and the emerging dominance of restricted stock.

### 5.2.1 Managerial Power

The "managerial power" approach begins with the self-interested executives envisioned by Berle and Means (1932) and Jensen and Meckling (1976) and adds a new element: the ability of these executives to influence both the level and composition of their own compensation packages, often (if not invariably) at the expense of shareholders. One of the early contributors to this view is David Yermack, who has argued that CEOs extract rents from shareholders by timing their option grants to occur just before the release of good news (Yermack, 1997), by insider trading through their family charitable foundations (Yermack, 2009), through lucrative severance and change in control provisions (Hartzell, Ofek, and Yermack, 2004; Yermack, 2006b; Dahiya and Yermack, 2008), and by consuming excessive perquisites (Yermack, 2006a).

The researchers most closely associated with the managerial power approach are Lucian Bebchuk and Jesse Fried, who have argued in a series of papers that both the level and composition of pay are determined not by competitive market forces but rather by captive board members catering to rent-seeking entrenched CEOs.[160] In addition, the authors argue that the CEO's ability to extract rent is limited by outside scrutiny and criticism (the "outrage constraint"), and CEOs respond by extracting rents

---

[160] See, for example, Bebchuk and Fried (2003, 2004a, 2004b), Bebchuk, Grinstein, and Peyer (2010), Bebchuk, Fried, and Walker (2002), and Fried (1998, 2008a, 2008b).

through difficult-to-observe or assess forms of compensation rather than through increased base salaries. They use their model to explain several common features of executive compensation plans, including the use (and misuse) of compensation consultants, the prevalence of stealth compensation (pensions, deferred pay, perquisites, and loans), gratuitous severance payments, and stock options that are uniformly granted at the money and not indexed for the market or industry.

*Can managerial power explain the trends in CEO pay?* There is no doubt that executives (like the rest of us) are self-interested and would prefer higher compensation to lower compensation. There is also little doubt that—while CEOs are never explicitly involved in setting their own pay (even those sitting on their own compensation committees)—CEOs have subtle ways of influencing the compensation committee and the pay-setting process.[161] However, as emphasized by Holmstrom and Kaplan (2003) and Frydman and Jenter (2010), there is no evidence that boards have become weaker or more captive over time. Indeed, every measure of board independence has improved since the mid-1980s. As discussed in Sections 5.1.1, the fraction of outside directors serving on the average board had increased to 80% by the end of the 1990s, and the CEO was the sole insider in nearly half of all firms. Since IRS Section 162(m) in 1993 (which required independence as a prerequisite for deductibility), most compensation committees have been fully independent. The 2003 NYSE listing requirements and 2010 Dodd–Frank Section 952 are appropriately characterized as tightening the definition from "independent" to "really independent" to "really, really independent", reflecting a mistaken belief that true independence can be measured by an objective standard applicable across all publicly traded companies without regard to the individual director. The increase in board independence during the 1990s should reduce managerial influence over pay, suggesting that the trends in CEO pay over the period were not driven by managerial power. In addition, the secular increase in disciplinarily firings of poorly performing CEOs (Kaplan and Minton, 2011; Huson, Parrino, and Starks, 2001) offers no evidence that boards are becoming systematically more passive over time.

Moreover, it is worth noting that many of the most generous and widely criticized option and severance payouts over the past two decades have been the direct result of formal employment agreements negotiated with external candidates, and not deals reached with powerful incumbents. Indeed, Murphy and Zábojník (2008) attribute the increase in executive pay to the increased prevalence of hiring CEOs from outside the firm. During the 1970s, under 15% of newly appointed CEOs were hired externally.

---

[161] For example, Murphy (1999) observes that while "outside board members approach their jobs with diligence, intelligence, and integrity…judgment calls tend systematically to favor the CEO. Faced with a range of market data on competitive pay levels, committees tend to on the high side. Faced with a choice between a sensible compensation plan and a slightly inferior plan favored by the CEO, the committee will defer to management. Similarly, faced with a discretionary choice on bonus-pool funding, the committee will tend to over- rather than under-fund."

By, the late 1990s, nearly a third of all CEO appointments came from outside of the firm, suggesting increasing competition in the managerial labor market. While the Murphy–Zábojník results (discussed in the next section) are often cited as evidence for the "efficient-contracting" approach, they are also consistent with directors systematically overpaying (and over-protecting) CEOs brought in from outside the firm.

In fact, compensation committees almost invariably pay "too much" for newly appointed CEOs, especially for those hired from outside the firm. Corporate directors seeking new CEOs from outside typically hire a professional search firm to identify qualified candidates for the position (Khurana, 2002a, 2002b). The pool of qualified candidates is narrowed through extensive research, background and reference checks, and interviews until a single individual is selected for the position. Negotiations over pay typically begin only after the favored candidate is identified and told that he or she is to be the new CEO. Indeed, many times negotiations are still on-going when the appointment is announced publicly. At this point the board is effectively locked in to the particular candidate CEO, which dramatically shifts the bargaining power to the seller (the candidate) rather than the buyer (the firm). This procedure is a reasonable way to identify top candidates when "price" is not an issue, but is clearly a recipe for systematically paying too much for managerial talent.

The tendency to pay too much and to pay it in the wrong way is exacerbated by potential CEOs who hire skilled contracting agents to negotiate on their behalf. In contrast, compensation committees rarely retain their own expert negotiators. The outcome is what one would expect in a game where there is such a clear mismatch: no matter how well intentioned, the typical compensation committee is no match against a professional negotiator, and overly generous pay packages become ubiquitous. But often the problem is worse: the incoming CEO (and his professional agent) negotiate not with the compensation committee but rather with the company's general counsel or head of human resources, knowing they will report to the CEO when the contracting is complete.

Overpaying newly hired CEOs is an agency problem caused by directors paying the new hires with shareholder money rather than their own. It is not, however, a "managerial power" problem, since the board is not captive and these are arms' length negotiations with a non-incumbent CEO candidate. The distinction is important because the policy prescriptions are different: the solution to overpaying new hires is to strengthen the negotiation process, while the solution to managerial power is to weaken the incumbent CEO's influence. More importantly, the "problem" of overpaying (and over-protecting) new hires may be small compared to the costs of selecting the wrong CEO.

In any case, hiring managerial talent from either inside or outside the firm is expensive, and the price of talent increased significantly during the 1990s and early 2000s. Kaplan and Rauh (2010) and Kaplan (2008), for example, present evidence that the increased pay for top executives is comparable to pay trends for top lawyers, investment

bankers, hedge-fund managers, venture capitalists, private-equity managers, and athletes. The rise in incomes for top talent in these disparate sectors—most with active and mobile labor markets—cannot plausibly be explained by managerial power. It seems unproductive to attribute gains in these other sectors to competitive market forces while inventing a different explanation for the rise in CEO pay. Indeed, the secular increase in external CEO appointments documented by Murphy and Zábojník (2008) suggests that the managerial labor market is becoming more rather than less competitive.

*Can managerial power explain the growth in the use of stock options?* Bebchuk, Fried, and Walker (2002) suggest that firms can "camouflage" excessive pay by substituting stock options for cash compensation, under the theory that such grants are difficult to value and are easy to hide in annual disclosures. Under disclosure rules effective before 1992, information on option grants was indeed difficult to obtain.[162] However, the center-piece of the sweeping new disclosure rules introduced in October 1992 focused on option grants, and two new tables were added to the proxy statements to describe the details of both the grant and the number and value of options held at the end of the year. Bebchuk, Fried, and Walker (2002) would predict that options grants would fall as the amount of information increased. However, option grants escalated (rather than fell) following the new rules.

Bebchuk and Grinstein (2005) attempt to provide a managerial power explanation for the 1990s increase in stock options as follows. First, they argue that the stock market boom weakened the outrage constraint, giving executives more latitude to increase their own pay. Second, they argue that increasing compensation in the form of options caused less outrage than increasing base salaries, not because of "camouflage" but because options offered the possibility of improved incentives. When the market declined in 2000–2002, the outrage constraint strengthened as investors became less forgiving of perceived managerial over-reaching, stemming the escalation in both pay and the use of stock options. Bebchuk and Grinstein (2005) use this framework to explain the correlation between CEO pay and general stock-price movements, as illustrated in Figure 18 in Section 3.7.5. Their framework would therefore also predict an increase in pay and options during the 2003–2007 bull market, and yet pay increases were modest and options were declining over this period. They could, of course, provide arguments for the existence of an "outrage constraint" for this period that would explain why pay levels moderated and options

---

[162] In September 1983, the SEC had reduced the amount of information companies needed to disclose on executive stock options. From 1978 to 1983, the "summary compensation table" in the proxy statement included not only cash compensation but also the number of new options granted and the increase in the intrinsic value of options held. Under the 1983 "paperwork reduction" rules, the summary compensation table included only cash compensation, the number of options granted was moved to later in the proxy, and information on outstanding options (and changes in the value of outstanding options) was eliminated. For details on the 1983 rules, see Hudson, "SEC Rules Allow Concerns to Curb Pay Disclosure: Companies Likely to Divulge Less on Executive Fees, Incentives, and Stock Options," *Wall Street Journal* (1983).

were replaced by restricted stock. This points to a basic problem with the Bebchuk and Grinstein (2005) explanation (and the managerial–power hypothesis more generally): there is no principled way to refute any trend in pay given the authors' flexible (and unmeasurable) definition of both the "outrage constraint" and its importance.

### 5.2.2 Efficient Contracting

The "efficient contracting" camp maintains that the observed level and composition of compensation reflects a competitive equilibrium in the market for managerial talent, and that incentives are structured to optimize firm value. The survey article by Edmans and Gabaix (2009) considers optimal-contracting explanations for the pay practices criticized under the managerial power camp, and the survey article by Frydman and Jenter (2010) discuss how these theories can predict increases in CEO pay over time.

Unlike the "managerial power" camp, the "efficient contracting" camp is not neatly characterized by a well-defined set of authors or articles. The modern executive compensation literature paralleled the emerging agency theory literature, and the majority of CEO pay papers written since the 1980s have been explicitly or implicitly based on agency or optimal-contracting theories. Indeed, the managerial power approach largely evolved as researchers—perhaps beginning with Jensen and Murphy (1990b) and Yermack (1995)—uncovered anomalies seemingly inconsistent with optimal contracts.

*Can efficient contracting explain the trends in CEO pay?* Beyond optimal incentive contracts, the efficient contracting approach includes market equilibrium models of managerial productivity, matching, and sorting that predict secular increases in CEO pay. For example, Murphy and Zábojník (2008) and Frydman (2007) offer general equilibrium models attributing the increase in executive pay to the increased prevalence of hiring CEOs from outside the firm. In particular, both papers attribute the trend toward outside hiring as reflecting gradual changes in the nature of the CEO job, modeled as a shift in the relative importance of general "managerial capital" (human capital specific to CEO positions) over firm-specific capital (reflecting skills, knowledge, contacts, and experience valuable only within the organization). The shift in the relative importance of general vs. firm-specific managerial capital leads to fewer promotions, more external hires, and an increase in equilibrium average wages for CEOs relative to the wages of lower-level workers. Ultimately, while it is plausible that the increased prevalence of outside hiring will increase average wages (if nothing else, employers must always pay a premium when hiring from outside compared to promoting from within), it is less plausible that the doubling of outside hiring from the 1970s to the 2000s could lead to such a huge increase in real CEO pay over this time period.

Alternatively, Gabaix and Landier (2008) build an equilibrium model in which the marginal product of managerial ability increases with firm size (so that it is optimal to assign the most talented managers to the largest firms). As shown by Rosen (1981) and Rosen (1982), such assortative matching produces equilibrium wages that are convex in

ability, such that small increases in ability can lead to large increases in wages (since the CEO is assigned to a larger firm). Gabaix and Landier (2008)'s key insight is that the wage of a CEO will depend not only on firm size, but also on the size distribution of all firms in the relevant market: as the average firm becomes larger, managerial marginal products increase and competition for scarce managerial talent will bid up compensation. In particular, they show that a shift in the size distribution of firms will lead to a proportional shift in compensation, and conclude that "the six-fold increase in CEO pay between 1980 and 2003 can be fully attributed to the six-fold increase in market capitalization of large US companies".

Gabaix and Landier (2008)'s results are consistent with the near-perfect correlation between CEO pay and general stock-price movements observed from 1980 to 2002 (see Figure 18 in Section 3.7.5). However (and similar to the critique of Bebchuk and Grinstein (2005) above), their results are not consistent with time trends in CEO pay and the stock market since 2002. In addition, while their insights on the size distribution are potentially important, their focus on market capitalization as the size measure is problematic since it conflates size, stock-price performance, and the vagaries of the market. Few would argue, for example, that Apple was really the largest firm in the world economy in 2012 (and yet their market value in early 2012 eclipsed that of Exxon-Mobil, PetroChina, and Royal Dutch Shell). Similarly, Volkswagen was not the second-largest firm on the planet for a couple of days in late October 2008 after its stock price increased by 350% over a two-day period (before tumbling by 60% over the following week).[163] While average CEO pay may have moved roughly proportionately with average market capitalization between 1980 to 2003, it far outpaced the growth in more traditional measures of size. For example, average revenues for the 500 largest US firms ranked by revenue grew only by 50% after inflation from 1980 to 2003, while average employment for the 500 largest US employers grew only by 19%.[164]

*Can efficient contracting explain the growth in the use of stock options?* The CEOs in most market-equilibrium models (including Murphy and Zábojník (2008), Frydman (2007), Gabaix and Landier (2008), and the informal model in Kaplan and Rauh (2010)) contribute only ability and not effort. Therefore, there is no role for incentives and thus no obvious reason why the increase in pay would come in the form of increased equity-based compensation (or, in particular, in stock options and why the preferred form of equity incentives would shift to restricted stock after 2002). To "explain" trends in CEO pay, it is not enough to predict increases in the level of pay, independent of

[163] Zuckerman, Strasburg and Esterl, "VW's 348% Two-Day Gain Is Pain For Hedge Funds," *Wall Street Journal* (2008).
[164] The Top 500 are for all US-based firms in Compustat. Using the same methodology, I find that the average market value (including debt and equity) for the 500 largest US firms grew by 300% between 1980 and 2003, substantially less than the 500% alleged by Gabaix and Landier (2008). I am unable to reconcile the difference.

dramatic changes in its composition. Indeed, as discussed above in Section 2.1.2, CEOs naturally demand a "risk premium" for accepting stock options in lieu of safer forms of compensation, and this risk premium will increase when the CEO is less diversified (i.e. when holding more shares of stock, or when the value of option portfolios increase relative to other wealth). Therefore, any increase in stock options will naturally be associated with an increase in total compensation, especially in a rising market. As shown in Figure 7 in Section 2.1.2, median risk-adjusted CEO pay actually fell from 1998 to 2001 (at least given the assumptions in the figure), even as the median unadjusted pay was exploding. In fact, the puzzle to be solved in Figure 7 is not why pay levels increased in the late 1990s (because they actually *declined* after adjusting for risk), but rather why risk-adjusted pay levels increased dramatically from 2002–2007, as companies replaced risky stock options with less risky restricted stock, without substantial declines in the grant-date fair-market value of pay.

Optimal-contracting theory (i.e. the subset of efficient contracting predicting that incentives are structured to optimize firm value) offers few explanations for the increase in option-based pay (i.e. increases in pay-performance sensitivities) in the 1990s. Consider, for example, the benchmark model where firm value is given by $y = a + \varepsilon$, where $a$ is executive effort, and $\varepsilon$ is (normally distributed) uncontrollable noise, $\varepsilon \approx N(0,\sigma^2)$. Moreover, suppose that managerial contracts take the simple linear form $w(x) = s + by$, where s is a fixed salary and $b$ is the sharing rate (or "pay-performance sensitivity"). Assuming that the executive has exponential utility, $U(x) = -e^{r(w-c(a))}$, where $r$ is the executive's absolute risk aversion and $c(a)$ is the convex disutility of effort, the optimal sharing rate is given by:[165]

$$b = \frac{1}{1 + r\sigma^2 c''}.$$

Traditional contracting theory therefore suggests a finite number of factors that might explain higher incentives among CEOs in the 1990s. First, perhaps CEOs became less risk or effort averse in the 1990s, but to my knowledge there is no theory or empirical work suggesting such declines in risk- or effort-aversion parameters. Second, perhaps CEO performance became estimated with less noise in the 1990s. While potentially consistent with the increase in director independence (if taken as a proxy for board monitoring), most measures of cash-flow or shareholder-return volatility increased rather than decreased over this time period.

---

[165] For similar (but more general) derivations of the optimal pay-performance sharing råte, see Lazear and Rosen (1981), Holmstrom and Milgrom (1991), Gibbons and Murphy (1992), and Milgrom and Roberts (1992).

Alternatively, suppose that firm value is given by $y = \theta a + \varepsilon$, where the primary source of uncertainty is variations in $\theta$ (i.e. managerial productivity assumed to be observed by the CEO but not by directors or shareholders). Zábojník (1996) and Prendergast (2002) show that optimal pay-performance sensitivites increase with the volatility of $\theta$ (incentives are more important when the CEO has private information about his or her marginal productivity). Again, to my knowledge there is no theory or empirical work suggesting that CEO marginal productivity became more volatile during the early 1990s.

Moreover, optimal-contract theories must explain not only the increase in equity-based compensation, but why that increase came almost entirely in the form of stock options as opposed to restricted shares or other equity-based instruments. Several papers have attempted, with only limited success, to provide theoretical justification for stock options. For example, traditional principle-agent models based on constant relative risk aversion and lognormally distributed stock prices (e.g., Hall and Murphy, 2002; Dittmann and Maug, 2007), suggest that—when salaries can be adjusted—contracts with restricted shares or options granted in-the-money are generally are superior to contracts with at-the-money options.[166]

Ultimately, the most compelling optimal-contracting explanation for the increase in equity-based compensation in the 1990s is that contracts were *suboptimal* before the 1990s, and got better. As explored above in Sections 3.5.6 and 3.6, year-to-year changes in executive pay in the 1970s largely reflected changes in company revenues (rather than performance), contributing to unproductive diversification, expansion and investment programs. The takeover and LBO market of the 1980s demonstrated vast potential for value creation in previously inefficient firms, leading academics, institutions, and shareholder advocates to demand that pay be more closely tied to shareholder performance. As emphasized by Holmstrom and Kaplan (2001), stock options allowed executives to share in the value created by internal restructurings that reduced excess capacity or reversed ill-advised diversification programs. The growing importance of shareholder activists and large institutional investors (Gompers and Metrick, 2001) increasingly pressured firms to tie pay more closely to stock-price performance. Stock options also became the currency of choice for high-tech start-ups, rich with ideas but (allegedly) short of cash or sources of capital. As a result, the popularity of options soared with the stock market during the 1990s, to the benefit of shareholders and executives alike. In fact, part of the increase in options during the 1990s plausibly reflects the fact that they seemed to be working: corporate boards and top managers began to associate option grants with successful company performance, especially during the high-tech

---

[166] Contracting models justifying the use of stock options rather than stock typically focus on optimal risk taking rather than (or in addition to) effort incentives (see, for example, Hirshleifer and Suh, 1992; Edmans and Gabaix, 2011).

and Internet boom of the late 1990s. Indeed, the increase in options coupled with the renewed focus on shareholder value creation may help explain the overall growth in stock market during this period.

This optimal-contracting explanation for stock options cannot, however, explain the magnitude of the explosion, and why it came in the form of options rather than stock. Consider, for example:

- The increase in stock options for top-level executives was associated with no discernable decrease in other forms of compensation (such as base salaries, bonuses, or benefits). To my knowledge, there is not an efficient contracting theory that predicts stock options to be added "for free" on top of what were presumably competitive compensation packages.

- Most contracting models would predict that the number of options granted would decline as stock prices increase, since the Black–Scholes cost of granting at-the-money options increases proportionally with the stock price. However, the number of options (as a fraction of outstanding common stock) increased rather than decreased during the 1990s, leading to a near-perfect correlation between average option grant-date values and stock-market indices between 1980 and 2002 (see Figure 18).

- Beginning in 2002, and especially since 2006, restricted stock has replaced stock options as the dominant form of equity-based compensation (and, indeed, is now the largest single component of compensation for the typical CEO in S&P 500 firms). To my knowledge, there is not an efficient-contracting theory that predicts this switch.

Even more difficult for the optimal-contracting camp is explaining why so many options were granted to so many employees well below the executive suite (see Figures 19 and 20). Since non-tradable stock options are an unusually inefficient method of conveying compensation (see Section 2.1.2), the incentive benefit from stock options must exceed the substantial difference between the company's opportunity cost of granting options and the "value" of those options from the perspective of risk-averse undiversified employees. While there may be efficient contracting justifications for granting options to top-level executives and other critical employees who can directly impact company stock prices (such as R&D scientists), there is (to my knowledge) no compelling incentive theory explaining option grants for rank-and-file employees.

Existing theories of broad-based option plans focus not on incentives but on other aspects of the employment relation. Oyer (2004), for example, argues that broad-based options may help satisfy participation constraints when reservation wages are correlated with the "market" and when it is costly to adjust other terms of employee compensation. Oyer and Schaefer (2005) and Bergman and Jenter (2007) argue that it might be optimal to grant options rather than cash when employees are irrationally optimistic about company prospects. Core and Guay (2001) argue that firms grant options to

lower-level employees as a substitute for cash compensation, and document a greater use of options for firms facing financing constraints.[167] Babenko, Lemmon, and Tserlukevich (2011) argue that financially constrained firms rely on cash inflows from employee option exercises to finance investments.

The common failing in the aforementioned theories of broad-based stock option plans is neither recognizing nor incorporating the substantial difference between the company's cost and the employee's value of non-tradable stock options. For example, Oyer (2004) offers no compelling argument or evidence that options are an efficient substitute for flexible employment terms (indeed, he largely ignores the efficiency cost of options, and assumes that contract adjustments are prohibitively costly)[168], and Core et al. (2001) and Babenko et al. (2011) implicitly hold but provide no theoretical or conceptual evidence for the implausible assumption that risk-averse undiversified employees are efficient sources of capital. Bergman and Jenter (2007) suggest that firms can reduce compensation costs by paying over-optimistic employees with (potentially overvalued) options, but provide no evidence that options are offered as a substitute for other forms of compensation. Indeed, all these models ignore the fact that most broad-based option plans were layered on top of existing compensation arrangements, and were not substitutes for cash compensation. The dominant option granters in the 1990s were not small cash-poor internet start-ups (where a compelling incentive-based rationale for broad-based options can be made), but rather large cash-rich giants such as Microsoft, Intel, Cisco, and Apple.

### 5.2.3 Perceived Cost

In a series of papers—admittedly garnering less traction than either the "managerial power" and "efficient contracting" approaches—I've suggested an alternative explanation for the growth of option-granting in the 1990s: decisions over options were made based on the "perceived cost" of options rather than on their economic (or "opportunity") cost.[169] When a company grants an option to an employee, it bears an economic cost equal to what an outside investor would pay for the option. But, it bears no outlay of cash, and (prior to the 2006 changes in accounting rules) bears no accounting charge. Moreover, when the option is exercised, the company (usually) issues a new share to the executive, and receives both the exercise price and (for non-qualified stock options) a tax deduction for the spread between the stock price and the exercise price. These factors make the "perceived cost" of an option much lower than the economic cost.

---

[167] In contrast, Ittner, Lambert, and Larcker (2003) find that companies with *greater* cash flows use options more extensively.

[168] Indeed, Oyer (2004) should predict that stock options are a particularly *ineffective* substitute for *downward* adjustments in employment terms (presumably firms face fewer short-run costs of adjusting in employees' favor).

[169] See, in particular, Murphy (2002, 2003) and Hall and Murphy (2003).

From the perspective of many boards and top executives who perceive options to be nearly costless—or indeed deny that options have value when granted—the only way they can quantify the options they award is by the number of options granted. During the 1990s, the focus on the quantity rather than the cost of options was further solidified by the institutions that monitor option plans. For example (see Section 3.7.5), SEC disclosure rules in place between 1992 and 2006 required companies to report only the number of, rather than the value of, options granted in the "Summary Compensation Table", the primary or most visible compensation table in the company's annual proxy statement. Similarly (see Section 3.7.6), under the pre–2003 NYSE listing requirement companies must obtain shareholder approval for the total number of options available to be granted, but not for the cost of options to be granted.[170] In addition, advisory firms (such as Institutional Shareholder Services) often base their shareholder voting recommendations primarily on the option "overhang" (that is, the number of options granted plus options remaining to be granted as a percent of total shares outstanding), and not on the opportunity cost of the proposed plan. Therefore, boards and top executives often implicitly admitted that the *number* of options granted imposes a cost on the company, while at the same time denying that these options have any real dollar cost to the company.

In addition, boards and top executives understand that options, when exercised, dilute the shareholdings of current equity holders. The number of options granted is included in fully diluted shares outstanding and therefore increased grants will decrease fully diluted earnings per share. Thus the negative consequences associated with these reductions in earnings per share also vary with the number of options granted, and not with the dollar-cost of the grants, and are consistent with the observed excessive focus on the number of options awarded and outstanding and not their dollar cost to the firm.

The perceived-cost view of stock options explains why options were granted in such large quantities to large numbers of executives and employees and also explains why the grant-date opportunity cost of options rose dramatically and subsequently declined with the stock market from 1980 to 2003 as shown in Figure 18 in Section 3.7.5. If boards focused only on the number of options granted, and the number of options granted stayed constant or varied positively with stock market performance, then the cost of the annual option grants would rise and fall in proportion to the changes in stock prices.

The perceived-cost view also explains why the relation between executive pay and the S&P 500 Index shown in Figure 18 weakened beginning in 2003. As discussed in

---

[170] In addition, as discussed in Section 3.7.6, under the pre-2003 listing requirements, companies did not need shareholder approval for options that would be issued broadly to executives and employees throughout the organization, but only for option grants that would be concentrated among the highest-level executives.

Section 3.8.4, while FAS 123R required firms to expense their options beginning in 2006, many firms began voluntarily expensing in early 2003. Expensing options brings the perceived cost of options more in line with their opportunity cost, and companies responded to the robust stock market from 2003 to 2007 by decreasing the number of options granted as stock prices increased (rather than increasing the quantity of options as happened from 1993 to 2001). Moreover, expensing brings the accounting treatment of options in line with the accounting treatment of restricted stock, explaining the shift from options towards restricted stock.

Finally, the perceived-cost view explains many prevalent features of stock options offered by the managerial-power camp as evidence for their position. For example, Bebchuk, Fried, and Walker (2002) cite the scarcity of "indexed options" (i.e., options where the exercise price adjusts over time to market- or industry-wide price indices) as evidence for the managerial power hypothesis. However, prior to the 2006 imposition of FAS 123R, indexed options were subject to an accounting charge while traditional options were not, increasing the relative perceived cost of indexed options. Similarly, Bebchuk, Fried, and Walker (2002) suggest that firms use uniform option terms (e.g., granting options "at the money") because diverging from normal practice by granting in-the-money options would spark outrage. Under the perceived-cost view, companies grant at-the-money options to avoid the accounting expense associated with in-the-money options. Indeed, the unsavory practice of "backdating" (in which firms granted in-the-money options but retroactively set the exercise date so the options appeared to be granted at the money; see Section 3.8.2) allowed firms to convey a given level of compensation without an accounting charge using fewer options than would be required without backdating. While the apparently common practice subsequently became "criminalized", many of the participants at the time viewed the practice as a minor accounting transgression that saved the shareholders a little dilution.

The perceived-cost view is readily acknowledged by practitioners and compensation consultants, but is usually denied or dismissed by financial economists because it implies systematic suboptimal decision-making by managers and a fixation on accounting numbers that defies economic logic. But executives often respond to accounting concerns in ways that seem irrational to economists. For example (as discussed in Section 3.7.4), the practice of repricing options following stock downturns virtually disappeared in December 1998 after an accounting expense was imposed on repriced options, illustrating how companies respond to accounting rules that have no affect on company cash flows. Similarly (as discussed in Section 3.8.4), firms accelerated the exercisability of existing options in advance of the implementation of FAS 123R in order to avoid an accounting charge for previously granted but unexercisable options; such acceleration hurts shareholders by reducing retention incentives and allowing executives to unwind their equity positions. As another example (only slightly beyond the executive

compensation arena), companies systematically scaled back retiree healthcare benefits after FASB required companies to record a current accounting charge for anticipated future medical costs.[171] The new accounting rule apparently increased the perceived cost of these benefits, putting them more in line with their actual economic cost, and as a result companies reduced benefit levels.

While the perceived-cost approach can explain why so many options were granted to so many people (because options were free, or at least cheap), it cannot explain why the explosion in grants started in earnest in the early 1990s: after all, the accounting and tax rules governing non-qualified stock options had been in place since 1972. In addition, before the May 1991 ruling that allowed stock acquired by exercising options to be immediately sold (see Sections 3.5.4 and 3.7.2), companies routinely granted Stock Appreciation Rights (typically subject to an accounting charge) rather than stock options (typically subject to no accounting charge), suggesting that the choice of equity-based incentives were not solely driven by accounting considerations.[172]

More fundamentally, the problem with the perceived-cost approach is that stock options are, of course, neither free nor even cheap to grant. Indeed, non-tradable options are an unusually expensive way to convey compensation to risk-averse and undiversified employees (Hall and Murphy, 2002; Section 2.1.2 above). A tempting theory—consistent with the managerial power approach—is that executives fully understood the opportunity cost of options but duped gullible directors into believing they were free. However, this explanation is inconsistent with the fact that 95% of options were granted below the CEO level: it seems implausible that the CEO would support such as huge transfer of wealth from shareholders to employees for a modest increase in his or her own compensation.

More plausible is the idea that executives and directors simply misunderstood the nature of opportunity costs. There is ample evidence that executives routinely ignore the opportunity cost of equity capital, leading firms to excess capacity and inefficient levels of inventories, cash and working capital. Indeed, the "Economic Value Added" programs that became popular in the 1990s were explicitly designed to teach managers about the opportunity cost of capital. If executives have a hard time grasping the opportunity cost of equity, they will have an even harder time grasping the opportunity cost of a derivative on that equity, especially when told that the "cost" is not the accounting cost but rather is estimated using a seemingly arcane theoretical formula. It is worth recalling that—while the Black–Scholes methodology was twenty years old by the early 1990s and was increasingly being used in academic research on executive compensation—it

---

[171] See Amir (1993), Espahbodie, Strock, and Tehranian (1991) and Mittelstaedt, Nichols, and Regier (1995) for descriptions and analyses of FAS 106 (*Employers' Accounting for Postretirement Benefits Other than Pensions*).

[172] In particular, the accounting expense for SARs reflected the appreciation in stock prices from the grant date through the exercise date.

had only recently gained limited traction among compensation consultants, and was not considered a useful tool in most corporate human resources departments.

### 5.2.4  Politics of Pay

A central theme in this study has been the futility of "explaining" CEO pay without explicit consideration of the causes and consequences of government intervention into executive compensation over the past century. The option explosion in the 1990s, which in turn caused the escalation in pay levels that spawned both the efficient contracting and managerial–power literatures, is a prime example of this futility. In Section 3.7 I discuss six factors that I believe contributed to the 1990s explosion in stock options (and hence the escalation in pay):

- *Shareholder pressure for equity-based pay.* The takeover and LBO market of the 1980s demonstrated vast potential for value creation in previously inefficient firms, leading academics, institutions, and shareholder advocates to demand that pay be more closely tied to shareholder performance.
- *SEC holding-period rules.* In 1991, the SEC determined that shares acquired by exercising options could be sold immediately upon exercise (effectively eliminating the six–month holding requirement).
- *SEC option disclosure rules.* In 1992, the SEC required disclosure of only the number of options granted, and not the value of options granted. The new rules pre-empted a popular Senate bill demanding a single dollar value for total compensation (which, in turn, required a dollar-valuation for options).
- *Clinton's $1 million deductibility cap.* In 1993, Section 162(m) (which ironically was imposed to reduce levels of executive pay) provided a safe harbor for stock options, by exempting options from the $ 1 million deductibility limit.
- *Accounting rules for options.* In 1995, after pushing for expensing the "fair-market value" of stock options, FASB backed down and allowed options to be granted without an accounting expense to the company (thus preserving the illusion that options were nearly costless to grant).
- *NYSE listing requirements.* Under listing rules in place during the 1990s, companies needed shareholder approval for equity plans covering top-level executives, but did not need approval as long as a sufficient percentage of eligible employees were non-executives. Therefore, companies could bypass shareholder votes by granting options to lower-level employees as well as executives.

The first of these factors ("shareholder pressure for equity–based pay") is consistent with the efficient contracting explanation (at least the version of the theory that contracts were suboptimal before the 1990s, and got better). However, the remaining factors reflect government intervention into the pay process, often as unintended consequences of attempts to curb perceived excesses in executive pay (and executive stock options in particular).

For example, the May 1991 SEC rules that allowed executives to sell shares immediately after exercising options was an unintended consequence of an attempt to curb excessive grants. As discussed in Section 3.7.2, corporate insiders are required to report stock purchases on SEC Form 4, but were not (before May 1991) required to report option grants. To provide more transparency for option grants, the SEC redefined the "stock purchase" as the date the option was granted rather than when it was exercised (thus triggering Form 4 disclosure of grants within 10 days of the end of the month when options were granted). As a result of this new definition, the six-month holding period required by the Securities Act started when the option was granted and not when it was exercised, allowing immediate sales upon exercise and greatly enhancing the appeal of options.

Similarly, Bill Clinton's campaign promise to limit deductibility of executive pay covered all forms of pay, and was only later modified to exempt deductibility limitations for pay tied to productivity. After substantial debate, stock options with an exercise price at or exceeding the grant-date market price were defined as related to productivity, while options with a lower exercise price were (arbitrarily) defined as non-performance related. But, as discussed in Section 3.7.3, the intent of the Congressional sponsors of the ultimate legislation was to reduce "excessive compensation", and not to promote the use of stock options.

However, the government faced an interesting political quandary: while it sought to curb perceived excesses in executive pay and options, it simultaneously sought to encourage firms to issue options to lower-level employees. For example, in its 1992 disclosure rules, the SEC required firms to report not only the number of options granted to each proxy-named executive, but also report that number as a percentage of options granted to all employees. The sole purpose of this requirement—similar to the Dodd–Frank requirement to report the ratio of CEO pay to the pay of the median employee—was to encourage (or "shame") companies into spreading awards more equally across the organization.

The NYSE listing requirements—which required shareholder approval for executive option plans but not broad-based option plans—were also designed to encourage option grants to lower-level employees. As discussed in Section 3.7.6, until January 1998 it had generally assumed that "broad-based plans" excluded substantial grants to top executives, which limited their use. The "clarifications" in 1998 (revised in 1999) defined how companies could grant top-executive options without approval, so long as a sufficient percentage of either the eligible employees or options granted were below the top-executive level. As a consequence, grants to both executives and lower-level employees escalated.

Similarly, FASB's 1995 compromise (which allowed companies to continue to grant options without an accounting expense, while recommending expensing fair market values) was driven primarily by concerns about expensing's implications for lower-level

grants (and not concerns with top executives). Countering Carl Levin's (D-MI) Corporate Pay Responsibility Act requiring option expensing (Section 3.7.4), bilis were introduced in both the House and Senate against expensing. In May 1994, the US Senate passed (by a 88-9 vote) a non-binding "sense of Congress" resolution demanding FASB to drop its expensing proposal, claiming that expensing would affect the ability of companies to raise capital, create jobs, and attract the best employees.[173] The Senate was joined by the Clinton administration—in no means an advocate of high CEO pay—concerned that FASB's proposal would hurt the competitiveness high-tech companies.[174]

The political obsession for broad-based option programs continued into the early 2000s, even as the popularity of options waned due to stock-market declines and pressures towards voluntary expensing (Section 3.8.4). Advocates of broad-based plans in Congress, fearing that fair-market-value accounting for options would end of option grants to low-level employees, introduced several (ultimately shelved) bills to protect such programs, including:

- The Workplace Employee Stock Option Act of 2002 (H.R. 5242), which provided incentives for broad-based option programs by allowing employees to purchase options and stock through pre-tax payroll deductions, and providing accelerated tax deductions for employers.
- The Rank-and-File Stock Option Act of 2002 (S. 2877), which limited the tax deduction companies could take if a stock-option program was not broad based.
  These bilis, and several others, were shelved in committee and the factors that had encouraged broad-based options were reversed:
- NYSE and NASDAQ listing rules revised in 2003 required shareholder approval for all option plans (including broad-based plans);
- The SEC's 2006 disclosure rules required disclosure of grant-date values (and dropped the disclosure of the option grants to top executives as a percentage of all option grants);
- FASB revised its accounting rules effective for most companies in fiscal 2006, mandating the expensing of options at their grant-date fair market value.

Ultimately and predictably, these changes curtailed the practice of broad-based option plans: firms that already had such plans granted fewer options, and virtually no firms without plans introduced one. Indeed, as evident from Figure 19 in Section 3.7.6 the average number of options granted by firms to all employees in the S&P 500 fell by half from 2001 to 2005 (from 2.6% of outstanding shares each year in 2001 to 1.3% in 2005).

---

[173] "US Senate backs resolution to remove option plan," *Reuters News* (1994).
[174] "Clinton Enters Debate Over How Companies Reckon Stock Options," *Wall Street Journal* (1993).

## 5.3  Explaining Executive Compensation: It's Complicated

My objective in writing this study is to provide "context" for both research in executive compensation and the ongoing debate over pay. Executive compensation has evolved over time in response to changes in both economic and political environments. Most recent analyses of executive compensation have focused on efficient contracting or managerial-power rationales for pay, while ignoring or downplaying the causes and consequences of disclosure requirements, tax policies, accounting rules, legislation, and the general political climate. A central theme of this study is that government intervention has been both a response to and a major driver of time trends in executive compensation over the past century, and that any explanation for pay that ignores political factors is critically incomplete.

As an important example, the growth in stock options in the 1990s spawned a major literature focused on explaining both cross-sectional and time-series trends in equity-based compensation for US CEOs. This literature has largely ignored the importance of political factors. However, the initial popularity of stock options was a direct result of government policies in the 1950s (Section 3.4), as was the explosion (and subsequent implosion) of options in the 1990s and 2000s, respectively (Sections 3.7 and 3.8.4). Similarly, the contrasting evolution of stock options for US CEOs and their foreign counterparts is largely explained by political rather than economic factors (Section 4.3).

Indeed, what makes CEO pay both interesting and complicated is the fact that the efficient contracting, managerial power, and political paradigms co-exist and interact. In introducing plans that tie pay more strongly to performance as demanded by shareholders, directors routinely agree to pay more than necessary to compensate for the increased risk. Self-interested CEOs seek employment protection through overly generous severance provisions; directors acquiesce believing that the probability of failure is low (and because it is not their money anyway). When compensation failures occur (such as those overly generous severance payments), Congress gets outraged, triggering disproportionate reforms with little regard for shareholders or value creation. In turn, companies and their executives respond by circumventing or adapting to the reforms, usually in ways that increase pay levels and produce other unintended (and typically unproductive) consequences.

## REFERENCES

$1,623,753 Grace's Bonus For 1929: Bethlehem president testifies at merger trial to receiving this amount (1930). *Wall Street Journal* (July 22).

Aboody, D., Barth, M. E., & Kasznik, R. (2004). Firms' voluntary recognition of stock–based compensation expense. *Journal of Accounting Research, 42*, 123–150.

Abowd, J., & Bognanno, M. (1995). International differences in executive and managerial compensation. In R. Freeman & L. Katz (Eds.), Differences and changes in wage structures. The University of Chicago Press.

Abowd, J. M., & Kaplan, D. S. (1999). Executive compensation: Six questions that need answering. *Journal of Economic Perspectives, 13*, 145–168.

Aggarwal, R., Erel, I., Ferreira, M., & Matos, P. (2011). Does governance travel around the world? Evidence from institutional investors. *Journal of Financial Economics, 100*, 154–181.

Agrawal, A., & Mandelker, G. (1987). Managerial incentives and corporate investment and financing decisions. *Journal of Finance, 42*, 823–837.

Ailing options: Stock market decline dulls allure of plans for company officials (1957). *Wall Street Journal* (October 21).

Alchian, A. A., & Demsetz, H. (1972). Production, information costs, and economic organization. *American Economic Review, 62*, 777–795.

Alinsky wins at the SEC (2010). *Wall Street Journal* (August 30).

Alpern, R. L., & Gail, M. (2001). Guide to change of control: Protecting companies and their executives. Executive Compensation Advisory Services.

Amir, E. (1993). The market valuation of accounting information: The case of post-retirement benefits other than pensions. *The Accounting Review, 68*, 703–724.

Andrews, E. L., & Vikas, B. (2009). Amid fury, US is set to curb executives' pay after bailouts. *New York Times* (February 4).

Armstrong, C. S., Ittner, C. D., & Larcker, D. F. (2012). Corporate governance, compensation consultants, and CEO pay levels. *Review of Accounting Studies, 17*, 322–351.

Babenko, I., Lemmon, M., & Tserlukevich, Y. (2011). Employee stock options and investment. *Journal of Finance, 66*, 981–1009.

Baker, G. P., & Hall, B. (2004). CEO incentives and firm size. *Journal of Labor Economics, 22*, 767–798.

Baker, G. P., Jensen, M. C., & Murphy, K. J. (1988). Compensation and incentives: Practice vs theory. *Journal of Finance, 43*, 593–616. <http://papers.ssrn.com/Abstract=94029>.

Baker, J. C. (1938). *Executive salaries and bonus plans*. McGraw Hill.

Barboza, D. (2002). Enron's many strands: Executive compensation. Enron paid some, not all, deferred compensation. *New York Times* (February 13).

Bartov, E., Givoly, D., & Hayn, C. (2002). The rewards to meeting or beating earnings expectations. *Journal of Accounting and Economics, 33*, 173–204. <http://papers.ssrn.com/paper=247435>.

Bebchuk, L. A., & Fried, J. M. (2004a). *Pay without performance: The unfulfilled promise of executive compensation*. Cambridge, MA: Harvard University Press.

Bebchuk, L. A., & Fried, J. M. (2004b). Stealth compensation via retirement benefits. *Berkeley Business Law Journal, 1*, 291–326.

Bebchuk, L. A., Grinstein, Y., & Peyer, U. (2010). Lucky CEOs and lucky directors. *Journal of Finance, 65*, 2363–2401.

Bebchuk, L. A., Fried, J. M., & Walker, D. I. (2002). Managerial power and rent extraction in the design of executive compensation. *University of Chicago Law Review, 69*, 751–846. <http://papers.ssrn.com/abstract=316590>.

Bebchuk, L. A., & Grinstein, Y. (2005). The growth of executive pay. *Oxford Review of Economic Policy, 21*, 283–303.

Bebchuk, L. A., & Fried, J. M. (2003). Executive compensation as an agency problem. *Journal of Economic Perspectives, 17*, 71–92.

Bender, M. (1975a). The executive's tax-free perks: The IRS looks harder at the array of extras. *New York Times* (November 30).

Bender, M. (1975b). Fringe benefits at the top: Shareholder ire focuses on loan systems. *New York Times* (April 13).

Bentsen opposes FASB on reporting stock options (1993). *Wall Street Journal* (April 7).

Bergman, N., & Jenter, D. (2007). Employee sentiment and stock option compensation. *Journal of Financial Economics, 84*

Bergstresser, D., & Philippon, T. (2006). CEO incentives and earnings management. *Journal of Financial Economics, 80*, 511–529.

Berle, A. A., & Means, G. C. (1932). *The modern corporation and private property*. New York: Macmillan Publishing Co.

Berman, D. K. (2010). The game: New frontier for the SEC: The clawback. *Wall Street Journal* (June 22)

Berton, L. (1992). Business chiefs try to derail proposal on stock options. *Wall Street Journal* (February 5)

Berton, L. (1994). Accounting rule-making board's proposal draws fire. *Wall Street Journal* (January 5)

Bettner, J. (1981). Incentive stock options get mixed reviews, despite the tax break they offer executives. *Wall Street Journal* (August 24)

Bhagat, S., & Bolton, B. (2011). Bank executive compensation and capital requirements reform.

Big earners cashing in now: Fearful of Clinton's tax plans, they rush to exercise their options. *San Francisco Chronicle* (December 29, 1992).

Bizjak, J. M., & Anderson, R. C. (2003). An empirical examination of the role of the CEO and the compensation committee in structuring executive pay. *Journal of Banking and Finance* <http://ssrn.com/abstract=220851>.

Black, F., & Scholes, M. S. (1973). The pricing of options and corporate liabilities. *Journal of Political Economy, 81*, 637–654.

Blumenthal, R. (1977). Misuse of corporate jets by executives is drawing more fire. *New York Times* (May 19).

Board's text on executive compensation (1971). *Wall Street Journal* (December 28).

Bonus figures given at trial: Six vice presidents of Bethlehem received $1,432,033 in 1929, 1930. *Wall Street Journal* (July 23).

Botero, J., Djankov, S., La Porta, R., Lopez-De-Silanes, F., & Shleifer, A. (2004). The regulation of labor. *Quarterly Journal of Economics, 119*, 1339–1382.

Bowe, C., & White, B. (2007). Record payback over options, *Financial Times* (December 7).

Bray, C. (2007). Former Comverse official receives prison term in options case. *Wall Street Journal* (May 11)

Bryan, S., Nash, R., & Patel, A. (2006). The structure of executive compensation: International evidence from 1996–2004.

Bryant, A. (1998). New rules on stock options by Big Board irk investors. *New York Times* (April 22).

Burns, N., & Kedia, S. (2006). The impact of performance-based compensation on misreporting. *Journal of Financial Economics, 79*, 35–67.

Business groups oppose Nixon control plan, intensify their efforts to abolish restraints (1974). *Wall Street Journal* (February 25).

Cadman, B., Carter, M. E., & Hillegeist, S. (2010). The incentives of compensation consultants and CEO pay. *Journal of Accounting and Economics, 49*, 263–280.

Cai, J., & Vijh, A. (2005). Executive stock and option valuation in a two state-variable framework. *Journal of Derivatives*, 9–27.

Calame, B. (1972). Executives' pay faces going-over by wage board. *Wall Street Journal* (April 24)

Carter, M. E., & Lynch, L. J. (2003). The consequences of the FASB's 1998 proposal on accounting for stock option repricing. *Journal of Accounting & Economics, 35*, 51–72.

Choudhary, P., Rajgopal, S., & Venkatachalam, M. (2009). Accelerated vesting of employee stock options in anticipation of FAS 123R. *Journal of Accounting Research, 47*, 105–146.

Chrysler chairman defends option plan, offers to discuss it with federal officials (1963). *Wall Street Journal* (December 23).

Chrysler officers got profit of $4.2 million on option stock in '63 (1964). *Wall Street Journal* (January 15).

Chrysler officers' sale of option stock could stir tax bill debate (1963). *Wall Street Journal* (December 18).

Clinton enters debate over how companies reckon stock options (1993). *Wall Street Journal* (December 23).

Congress and taxes: Specialists mull ways to close loopholes in present tax laws (1959). *Wall Street Journal* (January 7).

Connor, J. E. (1987). There's no accounting for realism at the FASB. *Wall Street Journal* (March 26)

Conyon, M. J., Core, J. E., & Guay, W. R. (2011). Are US CEOs paid more than UK CEOs? Inferences from risk-adjusted pay. *Review of Financial Studies., 24*, 402–438.

Conyon, M. J., Fernandes, N., Ferreira, M. A., Matos, P., & Murphy, K. J. (2013). The executive compensation controversy: A transatlantic analysis. In T. Boeri, C. Lucifora, & K. J. Murphy (Eds.), *Executive remuneration and employee performance-related pay: A transatlantic analysis*. Oxford University Press.

Conyon, M. J., & Murphy, K. J. (2000). The prince and the pauper? CEO pay in the United States and United Kingdom. *Economic Journal, 110*, F640–F671.

Conyon, M. J., & Schwalbach, J. (2000). Executive compensation: Evidence from the UK and Germany. *Long Range Planning, 33*, 504–526.

Core, J., & Guay, W. (2001). Stock option plans for non-executive employees. *Journal of Financial Economics, 61*

Core, J., & Guay, W. (2002). Estimating the value of employee stock option portfolios and their sensitivities to price and volatility. *Journal of Accounting Research, 40*, 613–630.

Crystal, G. (1991). *In search of excess: The overcompensation of American executives*. New York: W.W. Norton & Company.

Crystal, G. S. (1984). Manager's journal: Congress thinks it knows best about executive compensation. *Wall Street Journal* (July 30)

Crystal, G. S. (1988). The Wacky, Wacky World of CEO Pay (June 6).

Cuomo, A. M. (2009). No rhyme or reason: The heads I win, tails you lose bank bonus culture (July 30).

Cut high salaries or get no loans, is RFC warning (1933). *New York Times* (May 29).

DavisPolk (2010). Summary of the Dodd–Frank Wall Street Reform and Consumer Protection Act, Enacted into Law on July 21, 2010 (July 21).

De Angelis, D., & Grinstein, Y. (2011). Pay for the right performance.

DeFusco, R., Johnson, R., & Zorn, T. (1990). The effect of executive stock option plans on stockholders and bondholders. *Journal of Finance, 45*, 617–627.

Dittmann, I., & Maug, E. (2007). Lower salaries and no options? on the optimal structure of executive pay. *Journal of Finance, 62*, 303–343.

Dittmann, I., & Yu, K.C. (2011). How important are risk-taking incentives in executive compensation? <http://ssrn.com/abstract=1176192>.

Djankov, S., La Porta, R., Lopez-De-Silanes, F., & Shleifer, A. (2008). The law and economics of self-dealing. *Journal of Financial Economics, 88*, 430–465.

Dye, R. A. (1992). Relative performance evaluation and project selection. *Journal of Accounting Research, 30*, 27–52.

Eckhouse, J. (1987). Tech firms' study: Accounting rule attacked. *San Francisco Chronicle* (April 10).

Edelson, R., & Whisenant, S. (2009). A study of companies with abnormally favorable patterns of executive stock option grant timing.

Edmans, A., & Gabaix, X. (2011). Tractability in incentive contracting. *Review of Financial Studies, 24*, 2865–2894.

Edmans, A. (2012). How to fix executive compensation. *Wall Street Journal* (February 27).

Edmans, A., & Gabaix, X. (2009). Is CEO pay really inefficient? A survey of new optimal contracting theories. *European Financial Management*, 15–16.

Edmans, A., & Gabaix, X. (2011). The effect of risk on the CEO market. *Review of Financial Studies, 24*, 2822–2863.

Edmans, A., Gabaix, X., & Landier, A. (2009). A multiplicative model of optimal CEO incentives in market equilibrium. *Review of Financial Studies, 22*, 4881–4917.

Edmans, A., Gabaix, X., Sadzik, T., & Sannikov, Y. (2012). Dynamic CEO compensation, *Journal of Finance* 67, 1593–1637. <http://ssrn.com/abstract=1361797>.

Edmans, A., & Liu, Q. (2011). Inside debt. *Review of Finance, 15*, 75–102.

Efendi, J., Srivastava, A., & Swanson, E. P. (2007). Why do corporate managers misstate financial statements? The role of option compensation and other factors. *Journal of Financial Economics, 85*, 667–708. <http://papers.ssrn.com/abstract=547922>.

Egelko, B. (2010). 18 months for ex-Brocade CEO. *San Francisco Chronicle* (June 25).

Elia, C. J. (1967). Opting for options: Stock plans continue in widespread favor despite tax changes. *Wall Street Journal* (July 15).

Erickson, M., Hanlon, M., & Maydew, E. L. (2006). Is there a link between executive compensation and accounting fraud? *Journal of Accounting Research, 44*, 113–143.

Espahbodie, H., Strock, E., & Tehranian, H. (1991). Impact on equity prices of pronouncements related to nonpension postretirement benefits. *Journal of Accounting & Economics, 4*, 323–346.

Excerpts from carter message to congress on proposals to change tax system (1978). *New York Times* (January 22).

Fahlenbrach, R., & Stulz, R. M. (2011). Bank CEO incentives and the credit crisis. *Journal of Financial Economics, 99*, 11–26.

Federal bureau asks salaries of big companies' executives (1933). *Chicago Daily Tribune* (October 18).

Fernandes, N. (2008). EC: Board composition and firm performance. The role of independent board members. *Journal of Multinational Financial Management, 18*, 30–44.

Fernandes, N., Ferreira, M. A., Matos, P., & Murphy, K. J. (2013). Are US CEOs paid more? An international perspective. *Review of Financial Studies* (forthcoming).

Ferri, F., & Maber, D. (2010). Say on pay votes and CEO compensation: Evidence from the UK.

Ferri, F., & Sandino, T. (2009). The impact of shareholder activism on financial reporting and compensation: The case of employee stock options expensing. *The Accounting Review, 84*, 433–466.

Fischel, D. (1995). Payback: The conspiracy to destroy Michael Milken and his financial revolution (Harper Business).

Fisher, L. M. (1986). Option proposal criticized. *New York Times* (December 27).

Fitzpatrick, D., Scannell, K., & Bray, C. (2010). Rakoff backs BofA accord, unhappily. *Wall Street Journal* (Febraury 23)

Flanigan, J. (1996). It's time for all employees to get stock options. *Los Angeles Times* (April 21).

Forelle, C. (2006). How journal found options pattern. *Wall Street Journal* (May 22)

Forelle, C., & Bandler, J. (2006). Backdating probe widens as two quit Silicon Valley firm; Power Integrations Officials leave amid options scandal; 10 companies involved so far. *Wall Street Journal* (May 6).

Fracassi, C., & Tate, G. (2012). External networking and internal firm governance. *Journal of Finance, 67*, 153–194.

Freudenheim, M. (1993). Experts see tax curbs on executives' pay as more political than fiscal. *New York Times* (February 12).

Fried, J. M. (2008a). Hands-off options. *Vanderbilt Law Review*, 61.

Fried, J. M. (2008b). Option backdating and its implications. *Washington and Lee Law Review*, 65.

Fried, J. M. (1998). Reducing the profitability of corporate insider trading through pretrading disclosure. *Southern California Law Review, 71*, 303–392.

Frydman, C. (2007). Rising through the ranks: The evolution of the market for corporate executives, 1936–2003.

Frydman, C., & Jenter, D. (2010). CEO compensation. *Annual Review of Financial Economics, 2*, 75–102.

Frydman, C., & Saks, R. (2005). Historical trends in executive compensation, 1936–2003. Harvard University Working Paper.

Frydman, C., & Saks, R. E. (2010). Executive compensation: A new view from a long-term perspective, 1936–2005. *Review of Financial Studies, 23*, 2099–2138. <http://ssrn.com/abstract=972399>.

Fugitive mogul's rent coup, 2009 New York Post (August 26).

Gabaix, X., & Landier, A. (2008). Why has CEO pay increased so much?. *Quarterly Journal of Economics, 123*, 49–100.

Gibbons, R., & Murphy, K. J. (1990). Relative performance evaluation for chief executive officers. *Industrial and Labor Relations Review, 43*, 30S–51S.

Gibbons, R., & Murphy, K. J. (1992). Optimal incentive contracts in the presence of career concerns: Theory and evidence. *Journal of Political Economy, 100*, 468–505.

Gompers, P. A., & Metrick, A. (2001). Institutional investors and equity prices. *Quarterly Journal of Economics, 116*, 229–259. <http://ssrn.com/abstract=93660>.

Government moves to hold executives to 5.5% pay boosts (1973). *Wall Street Journal* (August 31).

Grant, P., Bandler, J., & Forelle, C. (2006). Cablevision gave backdated grant to dead official. *Wall Street Journal* (September 22)

Greenhouse, S. (1993). Deduction proposal is softened. *New York Times* (April 9).

Grossman, S. J., & Hart, O. D. (1983). An analysis of the principal-agent problem. *Econometrica, 51*, 7–45.

Guay, W. R. (1999). The sensitivity of CEO wealth to equity risk: an analysis of the magnitude and determinants. *Journal of Financial Economics, 53*, 43–71.

Gupta, U., & Berton, L. (1986). Start-up firms fear change in accounting. *Wall Street Journal* (June 23)

Hall, B. J., & Liebman, J. B. (1998). Are CEOs really paid like bureaucrats?. *Quarterly Journal of Economics, 113*, 653–691.

Hall, B. J., & Murphy, K. J. (2002). Stock options for undiversified executives. *Journal of Accounting and Economics, 33*, 3–42. <http://papers.ssrn.com/abstract_id=252805>.

Hall, B. J., & Murphy, K. J. (2003). The trouble with stock options. *Journal of Economic Perspectives, 17*, 49–70.

Harlan, C. (1994). High anxiety: Accounting proposal stirs unusual uproar in executive suites. *Wall Street Journal* (March 7).

Harlan, C., & Berton, L. (1992). Accounting firms, investors criticize proposal on executives' stock options. *Wall Street Journal* (February 19).

Hart, O. D. (1983). The market mechanism as an incentive scheme. *Bell Journal of Economics, 14*, 366–382.

Hartzell, J., Ofek, E., & Yermack, D. (2004). Whats in it for me? CEOs whose firms are acquired. *Review of Financial Studies, 17*, 37–61.

Hartzell, J., & Starks, L. (2003). Institutional investors and executive compensation. *Journal of Finance, 58*, 2351–2374.

Healy, P. M. (1985). The effect of bonus schemes on accounting decisions. *Journal of Accounting & Economics, 7*, 85–112.

Hechinger, J., & Bandier, J. (2006). In Sycamore suit, memo points to backdating claims. *Wall Street Journal* (July 12).

Henning, P. J. (2010). Behind the fade-out of options backdating cases. *New York Times* (April 30).

Heron, R. A., & Lie, E. (2006a). Does backdating explain the stock price pattern around executive stock option grants? *Journal of Financial Economics*. <http://papers.ssrn.com/abstract=877889>.

Heron, R. A., & Erik, L. (2006b). What fraction of stock option grants to top executives have been back-dated or manipulated. Unpublished working paper.

Higgins, A. (2007). The effect of compensation consultants: A study of market share and compensation, policy advice (October).

Hirshleifer, D., & Suh, R. (1992). Risk, managerial effort, and project choice. *Journal of Financial Intermediation*, 308–345.

Hill, G. W. (1931). Got Bonus of $1,200,000 Stock. *New York Times* (March 13).

Hite, G. L., & Long, M. S. (1982). Taxes and executive stock options. *Journal of Accounting and Economics, 4*, 3–14.

Holderness, C. G., & Sheehan, D. P. (1985). Raiders or saviors? The evidence of six controversial investors. *Journal of Financial Economics, 14*, 555–579.

Holmstrom, B. (1979). Moral hazard and observability. *Bell Journal of Economics, 10*, 74–91.

Holmstrom, B. (1982). Moral hazard in teams. *Bell Journal of Economics, 10*, 74–91.

Holmstrom, B. (1992). Contracts and the market for executives: Comment. In Wein, Lars, & Wijkander, Hans, (Eds.), *Contract Economics*. Blackwell Publishers.

Holmstrom, B., & Kaplan, S. N. (2001). Corporate governance and merger activity in the United States: Making sense of the 1980s and 1990s. *Journal of Economic Perspectives, 15*, 121–144.

Holmstrom, B., & Milgrom, P. (1987). Aggregation and linearity in the provision of intertemporal incentives. *Econometrica, 55*, 303–328.

Holmstrom, B., & Milgrom, P. (1991). Multitask principal-agent analyses: Incentive contracts asset ownership, and job design. *Journal of Law, Economics, and Organization, 7*, 24–52.

Holmstrom, B. R., & Kaplan, S. N. (2003). The state of US corporate governance. What's right and what's wrong?. *Journal of Applied Corporate Finance, 5*, 8–20. <http://ssrn.com/abstract=441100>.

Holzer, J. (2011a). Corporate news: Court deals blow to SEC, activists. *Wall Street Journal* (July 23).

Holzer, J. (2011b). A yes in say on pay. *Wall Street Journal* (July 8).

Horstmeyer, D. (2011). Beyond independence: CEO influence and the internal operations of the board.

Hot topic: Probing stock-options backdating (2006). *Wall Street Journal* (May 27).

House group hears conflicting views on stock option taxes (1959). *Wall Street Journal* (December 8).

House Unit Seen Favoring Curbs on Stock Options (1963). *Wall Street Journal* (February 25).

Hudson, R. L. (1983). SEC rules allow concerns to curb pay disclosure: Companies likely to divulge less on executive fees, incentives, and stock options. *Wall Street Journal* (September 23).

Hulse, C., & Herszenhorn, D. M. (2008). Bailout plan is set; House braces for tough vote. *New York Times* (September 29).

Hunt, A. R. (1971). Board agrees on tightening of standards on executive pay, increases topping 5.5%. *Wall Street Journal* (December 17).

Huson, M., Parrino, R., & Starks, L. (2001). Internal monitoring mechanisms and CEO turnover: A long term perspective. *Journal of Finance, 56*, 2265–2297.

Hyatt, J. C. (1975). No strings: Firms lure executives by promising bonuses not linked to profits. *Wall Street Journal* (December 24).

Ittner, C., Lambert, R. A., & Larcker, D. F. (2003). The structure and performance consequences of equity grants to employees of new economy firms. *Journal of Accounting and Economics, 34*, 89–127.

Jensen, M. C. (1972). Bonuses rise through loopholes. *New York Times* (January 9).

Jensen, M. C. (1978). Executives' use of perquisites draws scrutiny. *New York Times* (April 24).

Jensen, M. C. (1986a). Agency costs of free cash flow: Corporate finance and takeovers. *American Economic Review, 76*, 323–329. <http://ssrn.com/Abstract=99580>.

Jensen, M. C. (1986b). The takeover controversy: Analysis and evidence. In J. Coffee L. Lowenstein & S. Rose-Ackerman (Eds.), Takeovers and contests for corporate control. New York: Oxford University Press.

Jensen, M. C. (1993). The modern industrial revolution exit and the failure of internal control systems. *Journal of Finance, 6*, 831–880. <http://papers.ssrn.com/Abstract=93988>.

Jensen, M. C. (2003). Paying people to lie: The truth about the budgeting process. *European Financial Management, 9*, 379–406. <http://papers.ssrn.com/Abstract=267651>.

Jensen, M. C., & Meckling, W. H. (1976). Theory of the firm: Managerial behavior, agency costs, and ownership structure. *Journal of Financial Economics, 3*, 305–360. <http://papers.ssrn.com/Abstract=94043>.

Jensen, M. C., & Murphy, K. J. (1990a). CEO incentives: It's not how much you pay. *But How, Harvard Business Review, 68*, 138–153. <http://papers.ssrn.com/Abstract=146148>.

Jensen, M. C., & Murphy, K. J. (1990b). Performance pay and top management incentives. *Journal of Political Economy, 98*, 225–265. <http://papers.ssrn.com/Abstract=94009>.

Jensen, M. C., & Murphy, K. J. (2012). The earnings management game, Harvard business school working paper; USC Marshall School Working Paper. <http://ssrn.com/abstract=1894304>.

Johnson, S. A., Ryan, H. E., & Tian, Y. S. (2009). Managerial incentives and corporate fraud: The sources of incentives matter. *Review of Finance, 13*, 115–145. <http://ssrn.com/abstract=395960>.

Johnston, D. C. (1998). Fast deadline on options repricing: As of next Tuesday, it's ruled an expense. *New York Times* (December 8).

Joseph, R. (1978). US Industries faces queries on its perks at annual meeting. *Wall Street Journal* (April 20).

Kaplan, S. (1994). Top executive rewards and firm performance: A comparison of Japan and the US. *Journal of Political Economy, 102*, 510–546.

Kaplan, S. N. (2008). Are U.S. CEOs overpaid?. *Academy of management, perspectives*, 1–16.

Kaplan, S. N., & Minton, B. A. (2011). How has CEO turnover changed? *International Review of Finance*.

Kaplan, S. N., & Rauh, J. (2010). Wall street and main street: What contributes to the rise in the highest incomes?. *Review of Financial Studies, 23*, 1004–1050.

Kato, T., & Kim, W., & Lee, J.-H. (2006). Executive compensation and firm performance in Korea.

Kato, T., & Long, C. (2005). Executive compensation, firm performance, and corporate governance in China: Evidence from firms listed in the Shanghai and Shenzhen stock exchanges.

Kerr, S. (1975). On the folly of rewarding A, while hoping for B. *Academy of Management Journal, 18*, 769–783.

Khurana, R. (2002a). The curse of the superstar CEO. *Harvard Business Review*, 3–8.

Khurana, R. (2002b). *Searching for a corporate savior: The irrational quest for charismatic CEOs*. Princeton, NJ: Princeton University Press.

Korn, M. (2010). Diebold to pay $25 million penalty. *Wall Street Journal* (June 3).

La Porta, R., Lopez-De-Silanes, F., & Shleifer, A. (2006). What works in securities laws?. *Journal of Finance, 61*, 1–32.

La Porta, R., Lopez-De-Silanes, F., Shleifer, A., & Vishny, R. (1998). Law and finance. *Journal of Political Economy, 106*, 1113–1155.

Lambert, R. A., Larcker, D. F., & Verrecchia, R. E. (1991). Portfolio considerations in valuing executive compensation. *Journal of Accounting Research, 29*, 129–149.

Larcker, D. F., McCall, A. L., & Ormazabal, G. (2012). *The economic consequences of proxy advisor say-on-pay voting policies.*

Lazear, E. P. (1989). Pay equality and industrial politics. *Journal of Political Economy, 97*, 561–580.

Lazear, E. P., & Rosen, S. (1981). Rank-order tournaments as optimum labor contracts. *Journal of Political Economy, 89*, 841–864.

Lee, E. (2007). Option lawsuit give up details: Shareholders suing Mercury Interactive over timing of grants. *San Francisco Chronicle* (February 21).

Lewellen, W. G. (1968). *Executive compensation in large industrial companies*. New York: National Bureau of Economic Research.

Lie, E. (2005). On the timing of CEO stock options awards. *Management Science, 51*, 802–812.

Maremont, M. (2005). Authorities probe improper backdating of options: Practice allows executives to bolster their stock gains; a highly beneficial pattern. *Wall Street Journal* (November 11).

Maremont, M. (2009). Backdating likely more widespread. *Wall Street Journal* (August 18).

Martin, R. L. (2011). *Fixing the game: Bubbles, crashes, and what capitalism can learn from the NFL.* Harvard Business Review Press.

McGahran, K. (1988). SEC disclosure regulation and management perquisites. *Accounting Review, 63*, 23–41.

Merton, R. C. (1973). The theory of rational option pricing. *Bell Journal of Economics and Management Science, 4*, 141–183.

Metz, T. (1978). Close look expected at executive perks in proxy material: SEC stress on disclosure is linked to coming tales of holder-assisted goodies. *Wall Street Journal* (February 27).

Meulbroek, L. K. (2001). The efficiency of equity-linked compensation: Understanding the full cost of awarding executive stock options. *Financial Management*, 5–44.

Milgrom, P., & Roberts, J. (1992). *Economics, organization, and management*. Englewood Cliffs, NJ: Prentice-Hall.

Mirrlees, J. (1976). The optimal structure of incentives and authority within an organization. *The Bell Journal of Economics, 7*, 105–131.

Mittelstaedt, F., Nichols, W., & Regier, P. (1995). SFAS no. 106 and benefit reductions in employer-sponsored retiree health care plans. *The Accounting Review, 70*, 535–556.

Mullaney, T. E. (1951). Parley here indicates the continued spread in industry of stock purchase option plans. *New York Times* (August 12).

Murphy, K. J. (1985). Corporate performance and managerial remuneration: An empirical analysis. *Journal of Accounting and Economics, 7*, 11–42.

Murphy, K. J. (1996). Reporting choice and the 1992 proxy disclosure rules. *Journal of Accounting, Auditing, and Finance, 11*, 497–515.

Murphy, K. J. (1999). Executive compensation. In A. Orley, & C. David (Eds.), *Handbook of labor economics*. North Holland.

Murphy, K. J. (2000). Performance standards in incentive contracts. *Journal of Accounting & Economics, 30*, 245–278. <http://papers.ssrn.com/abstract=189808>.

Murphy, K. J. (2002). Explaining executive compensation: Managerial power vs. the perceived cost of stock options. *University of Chicago Law Review, 69*, 847–869.

Murphy, K. J. (2003). Stock-based pay in new economy firms. *Journal of Accounting & Economics, 34*, 129–147.

Murphy, K. J. (2012). *Pay, politics and the financial crisis*. In A. Blinder, A. Lo, & R. Solow (Eds.), *Economic lessons from the financial crisis*. Russell Sage Foundation.

Murphy, K. J., & Jensen, M. C. (2011). CEO bonus plans and how to fix them, Harvard business school NOM unit working paper 12-022; Marshall school of business working paper no. FBE 02-11. Available at SSRN: <http://ssrn.com/abstract=1935654>.

Murphy, K. J., & Oyer, P. (2004). Discretion in executive incentive contracts. USC working paper.

Murphy, K. J., & Sandino, T. (2010). Executive pay and independent compensation consultants. *Journal of Accounting and Economics, 49*, 247–262.

Murphy, K. J., & T. Sandino (2012). Are Compensation Consultants to Blame for High CEO pay?

Murphy, K. J., & Zábojník, J. (2008). Managerial capital and the market for CEOs.

Muslu, V. (2008). Inside board membership, pay disclosures and incentive compensation in Europe.

Narayanan, M. P., & Seyhun, H. N. (2005). Effect of Sarbanes-Oxley act on the influencing of executive compensation.

**Research Paper**

**ECONOMICS**

# A Literature Review of Agency Theory

| **Shubhi Agarwal** | RESEARCH SCHOLAR, DEPARTMENT OF ECONOMICS, CCS UNIVERSITY, MEERUT |
| --- | --- |
| **Rohit Goel** | ASSISTANT PROFESSOR, DESHBANDHU COLLEGE, UNIVERSITY OF DELHI |
| **Pushpendra Kumar Vashishtha** | ASSISTANT PROFESSOR, KAMLA NEHRU COLLEGE, UNIVERSITY OF DELHI |

**ABSTRACT**

Agency theory is a set of proposition in governing a modern corporation which is typically characterized by large number of shareholders or owners who allow separate individuals to control and direct the use of their collective capital for future gains. These individuals may not always own shares but may possess relevant professional skills in managing the corporation. The theory offers many useful ways to examine the relationship between owners and managers and verify how the final objective of maximizing the returns to the owners is achieved. This paper reviews the extensive literature of agency theory along with some of its limitations and it also focuses that a firm can improve its performance if agency cost may be reduced.

**INTRODUCTION-**

A firm can be owned by a single person or by more than one person. A firm owned by a single person is called a sole proprietorship concern. In this case the owner is the manager and his interests are no different from that of the firm i.e., maximizing the firm value. But in majority of cases a single individual cannot provide the entire capital, expertise and resources; and hence few individuals, with similar objective, collectively carry out the business. A large number of investor provides the risk capital. They are called shareholders. Shareholders have the residual claim on the assets of the company. Therefore, the right to control the use of the assets of the firm vests in them. They are deemed owners of the company. Shareholders delegate the power to manage the company to board of directors. The board delegates the same to managers while retaining its role to monitor and control the executive management. Corporate governance literature views shareholders as the principal and manager as their agent and describes the relationship as principal-agent relationship-

"An agency relationship is defined as one in which one or more persons (the principals) engage another person (the agent) to perform some service on their behalf which involves delegating some decision making authority to the agent. ( Hill and Jones, 1992)".

The divergence of interest between the owners and the managers, due to the separation of ownership from control, results in the agency costs.

Dealing with the agency problem is not free. Unfortunately, there is an agency cost associated with coping with the agency problem. Agency costs usually fall under the category of operating expenses. If employees of a company take a business trip and book themselves into the most expensive hotel they can find or if they insist on the best computer in the market for their offices, those are examples of agency costs. Those things don't maximize the wealth of the shareholders but instead minimize it.

**LITERATURE REVIEW OF AGENCY THEORY-** The agency problem inherent in the separation of ownership and control of assets was recognised as far back as in the 18th century

by Adam Smith in his Wealth of Nations, and studies such as those by Berle and Means (1934) and Lorsch and Maclver (1989) show the extent to which this separation has become manifest in firms throughout the world. Under this agency relationship, both the agents and the principals are assumed to be motivated solely by self-interest. As a result, when principal delegates some decision making responsibility to the agents, agents often use this power to promote their own well-being by choosing such actions which may or may not in the best interests of principals (Barnea, Haugen and Sanbet, 1985; Bromwich, 1992; Chowdhury, 2004). In agency relationship, the principals and agents are also assumed to be rational economic persons who are capable of forming unbiased expectations regarding the impact of agency problems together with the associated future value of their wealth (Barnea et al., 1985). Agency theory is concerned with the contractual relationship between two or more persons. Jensen and Meckling (1976, p.308) define agency relationship as a contract under which one or more person (principals) engage another person (the agent) to perform some service on their behalf which involves delegating some decision making authority to the agent. Jensen and Meckling identify managers as the agents, who are employed to work for maximizing the returns to the shareholders, who are the principals. Jensen and Meckling assume that as agents do not own the corporations resources, they may commit moral-hazards (such as shirking duties to enjoy leisure and hiding inefficiency to avoid loss of rewards) merely to enhance their own personal wealth at the cost of their principal. To minimize the potential for such agency problems, Jensen (1983) recognizes two important steps-

1-The principal-agent risk-bearing mechanism must be monitored through the nexus of organization and contracts. The first step, considered as the formal agency literature, examines how much of risks should each party assume in return for their respective gains. The principal must transfer some rights to the agent who, in turn, must accept to carry out the duties enshrined in the rights.

2- The second step, which Jensen (1983, p. 334) identifies as the positive agency theory clarifies how firms use contractual monitoring and bonding to bear upon the structure designed in the first step and derive potential solutions to the agency

problems. The inevitable loss of firm value that arises with the agency problems along with the costs of contractual monitoring and bonding are defined as agency costs,(Jensen and Meckling, 1976).

The principal-agent problem is also an essential element of the incomplete contracts view of the firm developed by Coase (1937), Jensen and Meckling (1976), Fama and Jensen (1983 a,b), Williamson (1975,1985), Aghion and Bolton (1992), and Hart (1995). This is because the principal-agent problem would not arise if it were possible to write a complete contracts. In this case, the investor and the manager would just sign a contract that specifies ex-ante what the manager does with the funds, how the returns are divided up, etc. In other words, investor could use a contract to perfectly align the interests and objectives of managers with their own. However, complete contracts are unfeasible, since it is impossible to foresee or describe all future contingencies. This incompleteness of contracts means that investors and managers will have to allocate residual control rights in some way, where residual control rights are the rights to make decisions in unforeseen circumstances or in circumstances not covered by the contract.

**Limitations of agency theory-** There are a number of limitations of agency theory (Eisenhardt 1989; Shleifer and Vishny 1997; Daily et al. 2003):

- Agency theory assumes complete contracts (i.e. contracts that cater for all possible contingencies such as ambiguities in language, inadvertence, unforeseen circumstances, disputes, etc). Bounded rationality does not allow for complete and efficient contracts. Information asymmetries, transaction costs and fraud are insurmountable obstacles to efficient contracting.
- Agency theory assumes that contracting can eliminate agency costs. The many imperfections in the market indicate that this assumption is not valid.
- Third party effects are not recognised. Third parties are those affected by the contract but who are not party to the contract. Many boards are conscious of third party effects and adopt social as well as financial responsibilities. Thus, whereas Maximum economic efficiency may (theoretically) be achieved under agency theory, it will not achieve maximum social welfare.
- Shareholders are assumed to be only interested in financial performance.
- Directors and management are assumed to owe their duty to shareholders. The law requires that duty to be owed to companies.
- Boards have a number of roles. Agency theory may be suitable for the monitoring-of-managers role of boards, but it does not explain the other roles of boards. Agency theory is not informative with respect to directors resources, services and strategy roles.
- Much of the corporate governance research is conceptualised as deterrents to managerial self-interest. Agency theory treats managers as opportunistic, motivated solely by self-interest. Many would argue that this theory does not capture those who are loyal to their firms.
- Agency theory does not take account of competence. Thus, if even incompetent managers are honest (or are made honest by board control) they will still be limited in their ability to meet shareholder objectives. It is not enough to incentivise people to get a task done; they must have the ability to carry out the task (Hillman and Dalziel, 2003).

**CONCLUSION-** As per the agency theory, due to the divergence of interests and objectives of managers and shareholders, one would expect the separation of ownership and control to have damaging effects on the performance of firms. Therefore, one way of overcoming this problem is through direct shareholder monitoring via concentrated ownership. The difficulty with dispersed ownership is that the incentives to monitor management are weak. Shareholders have an incentive to free-ride in the hope that other shareholders will do the monitoring. This is because the benefits from monitoring are shared with all shareholders, whereas, the full costs of monitoring are incurred by those who monitor. These free-rider problems do not arise with concentrated ownership, since the majority shareholder captures most of the benefits associated with his monitoring efforts. Several mechanisms can reduce agency problems. An obvious one is managerial shareholdings. In addition, concentration shareholdings by institutions or by block holders can increase managerial monitoring and so improve firm performance, as an outsider representation on corporate boards. The use of debt financing can improve performance by inducing monitoring by lenders. The labour market for managers can motivate managers to attend to their reputations among prospective employers and so improve performance. Finally the threat of displacement imposed by market for corporate control can create a powerful discipline on poorly performing managers.

To conclude it can be said, if agency costs may be reduced in the corporations, the firm performance can be improved.

## REFERENCES

• 599-616,Barnea, A., Haugen, R.A., and Senbet, L.W. (1985). Agency Problems and Financial Contracting, New Jersey: Prentice. | • Berle, A.A., and Means, G.C.,(1934), The Modern Corporation and Private Property, New York: The Macmillan Company, 178-197. | • Bromwich, M. (1992), Financial Reporting, Information and Capital Markets, London: Pitman | • Chowdhury, D. (2004), Incentives, Control and Development: Governance in Private and Public Sector with Special Reference to Bangladesh Dhaka: Viswavidyalay Prakashana Samstha. | • Daily,C.,Dalton,D.R.,andCannella, A.A. (2003) Corporate Governance: Decades of Dialogue and Data, The Academy of Management Review 28(3): 371-82. | • Eisenhardt, K. (1989) Agency theory: an assessment and review, Academy of Management Review 14: 5774 | • Hill, C. W. L., & Jones, T. M. 1992., Stakeholder-agency theory. Journal of Management Studies, 29: 131-154. | • Hillman, A.J. and Dalziel, T. (2003) Boards of Directors and Firm Performance: Integrating Agency and Resource Dependence Perspectives, Academy of Management Review 28(3): 383-96. | • Jensen, M. and W. Meckling (1976), Theory of the firm: managerial behaviour, agency costs and ownership structure, Journal of Financial Economics, 3 pp. 305-360. | • Jensen, M.C., (1983), Organisation Theory and Methodology the Accounting Review, V.LVIII, 2, pp 319-339. | • Maher Maria and Andersson Thomas, Corporate Governance: Effects on Firm Performance and Economic Growth, OECD, 1999. | • Shleifer Hall, Inc. Beaver, W.H. (1989)., A. and Vishny, R. W. (1997), A Survey of Corporate Governance, Journal of Finance 52(2)(June): 737-77.

Nicklaus, D. (2010). Scandal left both sides sullied: Backdating undermined confidence, but some good guys overreached, St. Louis Post-Dispatch (February 21).

Nixon halts push to retain some of phase 4 controls (1974). *Wall Street Journal* (April 5).

Old wage board exits: new unit to take over with reduced powers (1952). *Wall Street Journal* (July 30).

One in 6 companies gives stock options (1952). *New York Times* (June 30).

Options defended at salary hearing: Restricted stock plans called neither inflationary nor compensatory by 8 men (1951). *New York Times* (August 7).

Options on stocks scored at hearing: Majority of witnesses call it inflationary and unfair to small stock holders (1951). *New York Times* (August 9).

Options on the wane: Fewer firms plan sale of stock to executives at fixed exercise prices (1960). *Wall Street Journal* (December 6).

Ostroff, J. (1993). Clinton's economic plan hits taxes, payrolls and perks (February 18).

Oyer, P. (2004). Why do firms use incentives that have no incentive effects?. *Journal of Finance, 59*, 1619–1649.

Oyer, P., & Schaefer, S. (2005). Why do some firms give stock options to all employees: An empirical examination of alternative theories. *Journal of Financial Economics, 76*, 99–133.

Peers, A. (1991). Executives take advantage of new rules on selling shares bought with options. *Wall Street Journal* (June 19)

Penn, S. (1978). Ford Motor covered upkeep for elegant co-op of chairman: Questions arise on personal vs. business use of suite in posh New York hotel. *Wall Street Journal* (April 24).

Perry, T., & Zenner, M. (2001). Pay for Performance? government regulation and the structure of compensation contracts. *Journal of Financial Economics, 62*, 453–488.

Personal-use perks for top executives are termed income: SEC says valuable privileges will have to be reported as compensation by firms (1977). *Wall Street Journal* (August 22).

Plitch, P. (2006). Paydirt: Sarbanes-Oxley a pussycat on clawbacks. *Dow Jones Newswires* (June 9).

Politics and policy-campaign '92: From Quayle to Clinton, politicians are pouncing on the hot issue of top executive's hefty salaries (1992). *Wall Street Journal* (January 15).

Prendergast, C. (2002). The tenuous trade-off between risk and incentives. *Journal of Political Economy, 110*, 1071–1102.

President studies high salary curb: Tax power is urged as means of controlling stipends in big industries (1933). *New York Times* (October 23).

Railroad salary report: ICC asks Class 1 roads about jobs paying more than $10,000 a year (1932). *Wall Street Journal* (April 28).

Rankin, D. (1978). Incentives for business spending proposed in corporate package. *New York Times* (January 22)

RFC fixed pay limits: Cuts required to obtain loans (1933). *Los Angeles Times* (May 29).

Ricklefs, R. (1975). Sweetening the pot: Stock options allure fades, so firms seek different incentives. *Wall Street Journal* (May 27)

Ricklefs, R. (1977). Firms offer packages of long-term incentives as stock options go sour for some executives. *Wall Street Journal* (May 9)

Robbins, L. H. (1933). Inquiry into high salaries pressed by the government. *New York Times* (October 29)

Rose, N. L., & Wolfram, C. D. (2002). Regulating executive pay: Using the tax code to influence chief executive officer compensation. *Journal of Labor Economics, 20*, S138–S175.

Rosen, S. (1981). The economics of superstars. *American Economic Review, 71*, 845–858.

Rosen, S. (1982). Authority, control, and the distribution of earnings. *Bell Journal of Economics, 13*, 311–323.

Ross, S. A. (1973). The economic theory of agency: The principal's problems. *American Economic Review, 62*, 134–139.

Rudnitsky, H., & Green, R. (1985). Options are free, aren't they? Forbes (August 26).

Rules are issued on stock options (1951). *New York Times* (November 15).

Ryst, S. (2006). How to clean up a scandal. BusinessWeek.com (November 27).

Salary board urged to ban stock option plans until end of emergency (1951). *Wall Street Journal* (August 9).

Salary board's panel to study stock option in top executive pay (1951). *Wall Street Journal* (July 17).

Saly, P. J. (1994). Repricing executive stock options in a down market. *Journal of Accounting and Economics, 18*, 325–356.

Scannell, K., Rappaport, L., & Bravin, J. (2009). Judge tosses out bonus deal—SEC pact with BofA over Merrill is slammed; New York weighs charges against Lewis. *Wall Street Journal* (September 15)

Scharfstein, D. S. (1988). Product market competition and managerial slack. *Rand Journal of Economics, 19*, 147–155.

Scheck, J., & Stecklow, S. (2008). Brocade Ex-CEO gets 21 months in prison. *Wall Street Journal* (January 17)

Schellhardt, T. D. (1977). Perilous perks: Those business payoffs didn't all go abroad; bosses got some, too; IRS and SEC investigating loans and lush amenities provided for executitves; an eye on hunting lodges. *Wall Street Journal* (May 2).

Scipio, P. (1998). NYSE opens option loop hole. Investor Relations Business (May 11).

SEC exempts rights to stock appreciation from insider curbs (1976). *Wall Street Journal* (December 29).

SEC to push for data on pay of executives (1992). *Wall Street Journal* (January 21).

Senate unit votes to tighten rules on stock options (1964, January 15).

Shareholder groups cheer SEC's moves on disclosure of executive compensation (1992). *Wall Street Journal* (February 14).

Siconolfi, M. (1992). Wall Street is upset by Clinton's support on ending tax break for excessive pay. *Wall Street Journal* (October 21)

Skinner, D. J., & Sloan, R. G. (2002). Earnings surprises growth expectations, and stock returns or don't let an earnings torpedo sink your portfolio. *Review of Accounting Studies, 7*, 289–312.

Smith, A. (1776). *The Wealth of Nations* (Modern Library, Edited by Edwin Cannan, 1904. Reprint edition 1937, New York).

Solomon, D., & Paletta, D. (2008). US bailout plan calms markets, but struggle looms over details. *Wall Street Journal* (September 20)

Stanton, T. (1964). Cash comeback: Stock options begin to lose favor in wake of tax law revision. *Wall Street Journal* (August 10)

Stewart, G. B. (1991). *The quest for value: A guide for senior managers*. New York: Harper Business.

Stock options: Industry says salary board should keep its hands off employee plans (1951). *Wall Street Journal* (August 7).

Sundaram, R., & Yermack, D. (2007). Pay me later: Inside debt and its role in managerial compensation. *Journal of Finance, 62*, 1551–1588.

Sycamore Networks (2001). Q2 stock option grants issues. <http://online.wsj.com/public/resources/documents/sycamore_memo071206.pdf>.

Taibbi, M. (2011). Politics: OWS's beef: Wall street isn't winning—it's cheating. *Rolling Stone* (October 25)

Thomas, R. S. (2008). International executive pay: Current practices and future trends.

Tse, T. M. (2009). Shareholders say yes to executive pay plans; review tracks advisory votes at TARP firms. *Washington Post* (September 26)

US Senate backs resolution to remove option plan (1994). *Reuters News* (May 4).

US Steel guards data on salaries: Sends details confidentially to SEC head with request that they be kept secret (1935). *New York Times* (June 2).

Waxman, H. A. et. al. (2007). Executive pay: Conflicts of interest among compensation consultants (December).

Wells, H. (2010). No man can be worth $1,000,000 a year: The fight over executive compensation in 1930s America. *University of Richmond Law Review*, 44.

Wells, H. (2011). US executive compensation in historical perspective. In J. Hill, & R. S. Thomas (Eds.), *The research handbook on executive pay*. Edgar Elgar.

Wotapka, D. (2010). Former CEO at KB Home is convicted. *Wall Street Journal* (April 22)

Yermack, D. (1995). Do corporations award CEO stock options effectively?. *Journal of Financial Economics, 39*, 237–269.

Yermack, D. (1997). Good timing: CEO stock option awards and company news announcements. *Journal of Finance, 52*, 449–476. <http://papers.ssrn.com/abstract=8189>.

Yermack, D. (2006a). Flights of fancy: Corporate jets, CEO perquisites, and inferior shareholder returns. *Journal of Financial Economics*, 80.

Yermack, D. (2006b). Golden handshakes: Separation pay for retired and dismissed CEOs. *Journal of Accounting and Economics, 41*, 237–256.

Yermack, D. (2009). Deductio ad absurdum: CEOs donating their own stock to their own family foundations. *Journal of Financial Economics*, 94.

Zábojník, J. (1996). Pay–performance sensitivity and production uncertainty. *Economic Letters, 53*, 291–296.

Zhou, X. (2000). CEO pay, firm size, and corporate performance: Evidence from Canada. *Canadian Journal of Economics, 33*, 213–251.

Zimmerman, F. L. (1975). Washington word: Don't do as we do but do as we say: for bureaucrats, lawmakers, hard times aren't here; limousines and free trips. *Wall Street Journal* (February 7)

Zuckerman, G., Strasburg, J., & Esterl, M. (2008). VW's 348% two-day gain is pain for hedge funds. *Wall Street Journal* (October 29)

*Full Length Research Paper*

# Role of the agency theory in implementing management's control

**Mohammad Namazi**

Department of Accounting, College of Economics, Management and Social Science, Shiraz University, Iran.

The major purpose of this article is to analyze the role of the "Agency Theory" in implementing effective control mechanisms. In effect, various control paradigms under the following situations are elaborated: a) When the agent's control system is merely based on the output under the condition of uncertainty; b) When the control mechanism is based on the output, and information about agent's action or effort; c) When the agent's monitoring control is based on the output, agent's action and additional variables. It is concluded that agency theory posits different organizational, behavioral, economical and controlling roles, and it is a potent framework which can be extricated in promulgation of the management control systems. Furthermore, the implementation of a control mechanism depends on the amount and contents of the public and/ or private information that exist in the domain of the managerial accounting system. The disseminated information and the concurrent variables surrounding the agency relations are also vital elements in creating any control system. Finally, the optimal control mechanism under each preceding conditions are revealed.

**Key words:** Agency theory, management control, accounting information systems, information asymmetry.

## INTRODUCTION

According to the management accounting literature (Zimmerman, 2010; Kaplan and Atkinson, 2012; Horngren, et al., 2012) one of the significant functions and responsibilities of the managers is exerting control over the firms' operations and resources. The classic definition of management control posits control as "the process by which managers assure that resources are obtained and used effectively in the accomplishment of the organization's objectives", (Anthony, 1965: 17).

Today, however, the domain of control has been extended to consider not only the operations and strategic positions of the company, but also contemplating behavioral issues, and providing incentives to employees. In essence, a suitable management accounting system should dovetail the "planning" and "control" system in such a manner to provide information concerning accountability, and feedback information to ensure that the company adapts the internal and external organizational and environmental changes, the employees' behavior,

and measurement of the firm's variances from the actual operations. Figure 1 demonstrates the elements of this system.

Figure 1 indicates that the major elements of a suitable management control system really consist of two cycles: 1) the strategic planning cycle, and 2) the control cycle (Horngren et al., 2012). The strategic planning cycle encompasses establishing long-term and short-term strategic objectives, measures and targets, and related standards and budgets. The control cycle consists of the components which are illustrated in Figure 1. Hence, the control cycle is started based upon the designated strategic planning. The feedback is the central focus of the control system since it obtains information from the strategic planning process and particularly from the standards and budgets- element that makes it possible to compare actual results with standards and budgets and to determine concurrent variances. The final step of the control cycle leads to assimilating relevant information for

E-mail: mnamazi@rose.shirazu.ac.ir. Tel: +98 711 6460520. Fax: +98 711 6460520

638

**Figure 1.** Major elements of the management control system.

future planning and it is also referred to the feedback component (Zimmerman, 2010).

Hence, in implementing a successful control system, many issues are significant, particularly the following: (Outley, 2006: 49)

1) What are the "performance criteria" that will represent suitable performance?
2) What are the "relevant standards" of performance?
3) What are the rewards, and "behavioral issues" that will lead to the successful attainment of the targets and objectives?

In effect, the major inquiry is: How can efficient control mechanisms be designed and implemented in order to attain the goals and/or objectives of the firms? And at the same time provide incentives to firms' individuals to adapt actions that would lead to the attainment of the goal congruency?

The major purpose of this article is to demonstrate that the proceeding inquiry can be addressed effectively via the agency theory paradigms (Baiman, 1982; Lan and Heracleous, 2010). In effect, it investigates the role of the agency theory in devising optimal-incentive control systems. The contributions of this investigation are as

follows:

1) It provides a comprehensive explanation concerning the delicate performance criteria under the uncertainty condition.

2) It explicitly and mathematically demonstrates relevant components of the significant variables and standards of the optimal control systems.

3) It quantitatively extricates the significance of contemplating rewards and behavioral issues in designing optimal-incentive control frameworks.

4) It unambiguously delineates the importance of pecuniary as well as non-pecuniary factors in implementing efficient control systems.

The organization of this article is as follows: the methodology of the study; the basic agency theory model; management's control via the basic agency paradigms; findings of the study under various control mechanisms; discussion and concluding remarks.

## METHODOLOGY

The methodology of this study is based upon implementation of the "agency theory" framework. This theory, as will be explained later, has demonstrated a potent potential for establishing management control systems, because it is rich and mathematically considers various pecuniary and non-pecuniary items existing in the control systems. In addition, this theory is selected because various accounting scholars (Dikolli, 2001; Eldenburg and Krishnan, 2003; Kren and Tyson, 2009) have discussed its significance for establishing efficient management control devices. They also have shown that agency theory is able to explain the holistic as well as embedded effects of the control issues, and has an extensive ability to capture various control mechanisms under the condition of uncertainty. The latter aspect is particularly important, since most management control variables are usually uncontrollable and would happen under uncertain conditions. Furthermore, given accounting information systems, agency theory designates the exact managements' control variables precisely, and determines the optimal control elements which could be established under various control situations. It also considers the behavior and motivational issues in establishing control mechanisms. In sum, no other theories is as rich as the agency theory in explaining the reasons for developing managements' control systems, considering their elements, and how they can be effective established in various organizations.

## BASIC AGENCY THEORY

Basic agency paradigm was developed in the economics literature during 1960s and 1970s in order to determine the optimal amount of the risk- sharing among different individuals (Spence and Zeckhauser, 1971; Ross, 1973; Jensen and Meckling, 1976; Harris and Raviv 1976, 1978; Holmstrom, 1979).

However, gradually the domain of the agency theory was extended to the management area for determining the cooperation between various people with different goals in the organization, and attainment of the goal

congruency (Eisenhardt, 1989). In 1980s, agency theory was also appeared extensively in the managerial accounting realms to determine the optimal-incentive contracting among different individuals and establishing suitable accounting control mechanisms to monitor their behaviors and actions (Demski, 1980; Biaman, 1982; Namazi, 1985). It is this last function of the agency theory that will be emphasized in this study.

In its primitive form, agency theory relates to situations in which one individual (called the agent) is engaged by another individual (called the principal) to act on his/her behalf based upon a designated fee schedule. Since both individuals are assumed to be utility maximizer, and motivated by pecuniary and non-pecuniary items, incentive problems may arise, particularly under the condition of uncertainty and informational asymmetry. That is, the objective function of the principal and the agent may be incompatible, and therefore, the agent may take actions which will jeopardize the principal's benefits.

In addition, an agency operates under the condition of risk and uncertainty. In effect, the basic agency theory usually assumes that both individuals are risk averse. Under this circumstances, the amount and content of the produced accounting information and other information sources would become a significant issue in risk sharing and controlling the agent's actions (Namazi, 1985; Baiman, 1982, 1990).

The preceding basic agency model, however, has also been extended to cases in which there are multiple agents (Holmstrom, 1979; Antle, 1982; Radner, 1981), private information (Penno, 1984), multiple period performance (Radner, 1981), and multi-objective models (Namazi, 1983). In addition, the effect of various cultures on the assumptions of the agency theory has also been investigated (Osterman, 2006; Kren and Tyson, 2009).

Given the agency theory paradigm, and following Alchian and Demsetz (1972), Jensen and Meckling (1972), and Kaplan (1984), among others, a firm can be characterized as a nexuses of contractual agreements among different individuals. In this view, contracts are considered as an appropriate means for resource allocation and revealing the scope of the firm's activities. In addition, they can be expended as a powerful framework for effective management accounting control mechanisms. In this context, performance measures, appropriate control variables, and exogenous and endogenous parameters affecting the control process, can be captured and specified quantitatively by adapting the "agency theory" framework. Hence, this study draws in the agency paradigms to investigate the role of the agency theory in establishing management's control.

## MANAGEMENTS' CONTROL VIA THE BASIC AGENCY PARADIGM

To illustrate the basic elements of a control model in

terms of the agency framework, consider a simple firm which consists of two individuals only; one individual is the owner (called principal), the other one is the manager (called the agent) . The principal has invested in the firm, and has delegated the responsibility of the decision making to the agent. The agent exerts his/her services based upon pre-specified contractual agreements. Since the agent is motivated primarily by his/her self interest, he may select actions which would jeopardize the principal's benefits. To prevent the agent, a suitable control mechanism must be established. Thus, agency theory provides a potent reason as to why maintaining control mechanisms is necessary. We define this role as "the organizational role" of the agency.

Let us assume the final outcome (revenues, profits, etc.) resulting from the firm's operations is x ε X. The manager's share is $x_m$, and the residual is paid to the owner (i-e., $x_0 = x - x_m$). The outcome is a function of: 1) the invested capital (q ε Q), 2) the manager's action (a ε A), 3) the manager's effort, (e ε E); and 4) the states of natures (s ε S). Thus, it can be represented as followas:

$$x = p \, (s, a, e, q) \tag{1}$$

The output (x) is reported by a designated accounting system or other information sources to both individuals- the principal and the agent- at some cost of reporting, c (r). To be effective, both the designated information system and it's generated signals must be efficient-that is, it should generate quality information which would reduce the amount of uncertainty and would entail the optimal risk sharing. We define this function as" the informational efficiency" role of the agency theory. Consequently, we have:

$$x = p \, (s, a, e, q) - c \, (r) \tag{2}$$

Informational asymmetry problems exist between the agent and the principal. Let n ε N and m ε M denote the information system which is possessed or accessed by the manager and owner respectively. The signals y ε Y (Y ε N) is obtained after the agent's actions and effort, but before $x_m$ is determined.

For risk sharing purposes, and devising an efficient control mechanism, the contract should be based on the "observable elements" by both individuals (Demski and Dye, 1999; Pacharn, 2008). These elements are "performance measures" that will be exerted to monitor the agent's output. We define this role as the "system evaluation" role of the agency theory. To illustrate this role mathematically, let $\hat{n}$ denotes the common information in the two systems and $\hat{y}$ the respected common signal. The management function can be represented as $\hat{y} = \hat{n}$ (s, a, e, q, x) and f = r (s, a, e, q, x). Consequently, contractual agreement can be characterized as

$x_m = r(\hat{y}, f)$. Since it is assumed that the outcome (x) is always observed by both individuals, x is included in $\hat{y}$.

Following the agency theory, it is also assumed that both the agent and principal are utility maximizer under Von Neumann-Morgenstern's utility axioms (Demski, 1980; Lan and Heracleous, 2010). The principle is concerned only with his/her net residuals, $x_0$, and is risk-neutral. However, the agent's utility is centered on his/her pecuniary return, $x_M$, as well as his/her act(a) and effort(e) exerted. It also can be represented as additive and differentiable equation as $U_M(x_M, a) = U(x_M) - V(a)$. The efforts expended by the agent are assumed to generate disutility to him/her. Incentive problems may arise here because the objectives of the agent and principal can be incompatible. Consequently, the agent might select actions which are more consistent with his/her self-interest, and less consistent with the principal's goals. The agent is assumed to be generally risk and work-averse. We define this role as "behavioral role" of the agency theory.

Under this condition, the expected utility of the manager who exerts effort (e) and action (a), can be represented as follows:

$$E(U_{Mh} \mid a, e, r) = \int_s U_M(r(n(a, e, q, h, s)), e) \phi s(ds) \tag{3}$$

If $e_h^*(I)$ represents the management's effort given the contract I and effort h, the expected utility of the managers is as follows:

$$E(U_{Mh} \mid e_M^*(I), I = \max_{a \in A} E(U_{Mh} \mid e, I)$$
$$(e \in E) \tag{4}$$

On the other hand, the expected utility of the owner is the function of:

1) the contract that he/she has signed with management, and 2) his/her estimation concerning the agent's effort during the contract period. If $e_h^*(I)$ represents the owner's estimation of the agent's effort, then the expected utility function of the owner is represented as follows:

$$E(U_0 \mid e^*(I), I) = \iint_{Hs} U_0(P(s, e_h^*(I), q, h)$$
$$-r(n(s, e_h^*(I), q, h))$$
$$-c(s, e_h^*(I), q, h, n))$$
$$\phi s(ds) \phi_0(h \mid I) dh \tag{5}$$

By solving Equations 4, and 5 simultaneously the Pareto-

**Table 1.** Illustration of the management-owner relations in a basic agency model.

| 1 | 2 | 3 |
|---|---|---|
| Management and owner select a contract and information system(s) cooperatively | Management selects action (a) independently | Management and owner jointly observe the outcome, and management's pay is based upon a designed contract at time 1 |

optimal contract[1] can be attained. This type of contracts not only maximizes the utility function of the agent and the principal, but also would lead to efficient allocations of the firm's resources. We define this role, as the "allocation role" of the agency theory.

Much of the early research on agency theory (Demski, 1980; Shavell, 1979, Holmstram, 1979) have attempted to solve the preceding problems. As a result, different cases have been investigated, which provide the optimal control mechanism for each situation. We define this role as "the optimal control monitoring selection" role of the agency theory. This role is explained next.

**Case 1: Management's control based on the output**

The primitive scenario of management's control occurs when only output (x) of the agency is expended for monitoring control. This situation is illustrated in Table1.

At time 1, the owner hires a manager in order to obtain the following benefits: 1) to increase the output and, therefore revenues; 2) to acquire management's knowledge, expertise, or information; 3) to obtain such non-pecuniary benefits as affiliation and prestige. Management, on the other hand, maintains a desire to acquire remuneration, a wish to transfer more of the risk to the owner, and a need to satisfy non- pecuniary items.

Management knows how his/her information, knowledge, and skills match the job to be performed. However, the owner cannot appraise these attributes directly with its control monitoring devices (including an accounting information system). If management is hired, at time 2, he/she can independently choose any action or effort level that maximizes his/her utility function. Although the owner cannot directly observe the management's action and effort levels at time 2, the payoff (or outcome) always becomes observable by both - management and owner- at time 3[2]. This outcome is a function of management's actions, efforts, and exogenous stochastic variables. The outcome may be revenues, profits, or products and services.

At time 1, the owner and management must also address two issues cooperatively: the kind of information

system(s) to be exerted for characterizing management's control mechanism, and the type of contract to be employed for (a) effecting an efficient risk sharing and (b) providing management with an appropriate incentive. This type of contract is known as a "Pareto-optimal incentive contract" (Namazi, 1985; Pacharn, 2008; Demski et al., 2009).

Mathematically, when an accounting system *ex ante* reports the value of (x), management's fee schedule for controlling mechanism must depend only on x. This dependency, however, interrupts any owner-management risk-sharing arrangements. If, for instance, management is risk-averse and the owner is risk-neutral, the optimal contract based on the output maximizes management's expected utility function without decreasing its welfare and, at the same time, provides management an incentive to take no action that jeopardizes the owner's well being.

Much of the research on agency theory has focused on solving the preceding obstacle, and a plausible solution has emerged. This solution, derived by Spence and Zeckhauser (1971) and Ross (1973), among others, assumes that owner and management have identical attitudes toward risk and that (x) is formally based on (a) and (s). Given the distribution of (s) and using the first-order condition, one can then determine an optimal fee schedule by applying the calculus of variations to a manager's specified program.

Wilson (1968), Ross (1973), Harris and Raviv (1976, 1978), demonstrate that under these circumstances, the Pareto-optimal fee schedule is linear, consisting of a fixed salary and a variable that designates profit sharing between the owner and management.

When only output is monitored for controlling mechanism, however, the owner cannot assert the effort level of management directly via the accounting information system. Thus, if a poor outcome results, management can always blame it on unfavorable states of nature rather than on lower effort. In agency theory, this unobservability of labor's effort (action) and its effect on the outcome is known as the "moral hazard" (Holmstrom, 1979; Demski and Dye, 1999). A fee schedule that subsumes moral hazard problems, is known as the "second-best solution" (Holmstrom, 1979; Shavell, 1979; and Demski, 1980) because it is based on an imperfect estimate of the agent's effort, and it compensates for the lack of observability of the agent's effort by trading some risk-sharing benefits to stimulate a proper level of effort. This is one disadvantage of a linear fee schedule.

[1] - A Pareto-optimal contract is a contract in which maximizes the expected utility functions of the both individuals - the agent and the principal- and at the same time, provides enough incentives for the agent not to take any actions that jeopardize the principal's benefits.
[2] - This assumption rules out the possibility of initiating contracts that are incompatible with incentive.

Demski (1980), however, demonstrates that when management's stewardship function is only based on the outcome (x), if management and owner are both risk neutral, and management possesses sufficient resources to acquire the firm by compensating the owner, establishing a "takeover contract"- which sells the whole firm to the agent and thus leaves no risk sharing- provides the "first-best solution". Consequently, he provides the following lemma:

**Lemma 1:** *In the basic stewardship model with the manager risk neutral, an appropriate takeover contract will produce the first-best solution (p. 107).*

He also maintains that if the manager is work-neutral (V(a) = constant) and is ready to implement the owner's fee schedule based on a fix amount r(0)= $k^x$, the "salary contract" is appropriate. Thus,

**Lemma 2:** *In the basic stewardship model with the manager work-neutral, an appropriate salary contract will produce the first-best solution (p. 107).*

These situations, in effect, provide the necessary conditions for designing effective management control systems.

## Case 2: Management's control based on the output and agent's action

In order to alleviate the consequences of the moral hazard, and to improve the performance criterion of the management's control system, the owner and management can select one of two finer accounting systems ($n_2$) that reports both the output (x) and management's action (a): (1) an accounting information system that produces an additional signal ($\hat{y}$), which reports a complete observation of the management's effort, or (2) one that extends the amount of accounting data to situations in which incomplete information ($\hat{y}$) about management's action is reported.

In the first case, ($n_2$) reports both (x) and the perfect value of (a). Consequently, r ($x, \hat{y}$) can be characterized by establishing a "forcing contract" based on the management's effort (action) level. If management supplies a predetermined effort level observable by the owner, he/she would receive a share based only on the outcome produced. On the other hand, if management fails to exert a proper level of effort, it would get nothing (Harris and Raviv, 1978: 24). Such a fee schedule, which entails optimal risk sharing between the individuals, becomes a "first-best solution" (Shavell, 1979; Demski, 1980; Stevens and Thevarajan, 2010). Demski (1980) has termed this as a "wage contract" and has provided the following lemma:

**Lemma 3:** *In the basic stewardship model with the effort and outcome jointly observed, an appropriate wage contract will produce the first-best solution (p.107).*

In addition, when perfect information about the states of nature exists, and the output (x) is observable by both parties, establishing a contract which guarantees a fixed income ($k^*$) for the agent would be optimal, as long as he/she exerts his/her best action. In this case, the optimal control mechanism is represented by

$$r(s, x) = x - p(s, \hat{a}, \overline{q}) + k^*.$$

Demski (1980) has termed this fee schedule as the "insurance contract" and provides the following lemma:

**Lemma 4:** *In the basic stewardship model with the state and outcome jointly observed, an appropriate insurance contract will produce the first-best solution (p. 107).*

Much of the earlier work in agency theory (Arrow, 1971; Stieglitz, 1975; Shavell, 1979) focused on determining the characteristics of the first-best solution. These characteristics will not, however, be discussed here since they cannot be applied unless management effort is known with complete certainty. As Holmstrom (1979) notes, in real situations, full information about labor's effort is either impossible or prohibitively expensive to obtain.

The fineness of an accounting information system, which is exerted as a potent control mechanism, however, can be enhanced by collecting additional imperfect information about management's effort and skills. Consequently, the owner's endeavor should be reducing the extent of the "moral hazard" and "adverse selection effects". Moral hazard occurs when principal cannot determine the exerted level of the agent's effort. Adverse selection effects is created when the agent claims his/her skill and experience is higher than the actual one, and the principal cannot ascertain and verify if the agent actually has expended his/her ability in the designated job (Demski and Feltham, 1978; Holmstrom, 1979; Demski, et al., 1999). This goal can he achieved by appropriating resources to generate or improve cost effective reports relating to the agent's effort. Since the owner cannot observe management's effort directly, it can conveniently establish management's standard of output and exert this standard to measure ex post production efficiency. This standard or budget can be utilized as a powerful control mechanism. Furthermore, the owner can instruct management to select an optimal level of action. This consideration allows the principal and agent to share the random fluctuations in output optimally. Subsequently, management can be paid a constant amount as a reward for selecting an optimal action. In addition, since the agent is aware that his effort is observable, and his payoff is therefore subject to less random noise, he/she will expend a level of effort that provides the owner with an expected utility higher than when effort is not observable (Diamond and Verrecchia, 1982; Kren and Tyson, 2009).

Observing management's effort, however, is not costless, and perhaps more important, not all additional

information has positive effects. Infact, additional incomplete information may provide an inaccurate signal concerning management's true effort (Shavell, 1979). Hence, an interesting inquiry arises: Given the new risk, can improvements (in the sense of the Pareto-optimal condition) be made by obtaining additional incomplete information about the management's effort? The answer relies on the value of the additional information.

Holmstrom (1979), for instance, has introduced the "In formativeness principle" to reply to this inquiry. This principle generally states that any information concerning the agent's selected action, effort and skill, although incomplete, should be included in the compensation contracts, as long as the marginal benefits of the information is greater than its marginal costs. In sum, agency theory (Demski and Feltham, 1978; Gjedsdal, 1981, Namazi, 1985, Baiman, 2000) supports collecting additional information about the agent's efforts and skills, because it will improve the fee schedule solution, and will also enhance the precision of the measurement criteria of the management's control systems.

## Case 3: Management's control based on the output, agent's action, and additional variables

This aspect extends the management-owner function to situations in which an accounting information system $(n_3)$ provides not only the information surrounding the output (x) and action (a), but also an additional signal $(y^*)$ that conveys information about the capital (q) provided by the owner. Assuming that owner has completely delegated decision-making related to the production process to management, and that the amount of (q) will not be altered throughout the production process, the issue becomes how a Pareto-optimal payment schedule between the owner and management can be devised, and given additional incomplete accounting information, what kind of contract a manager should select?

Demski and Feltham (1978) attempted to address these issues by formulating an agency model within a general framework at the aggregate level. They explicitly distinguished the agent's action (a) from the effort (e) and denoted the outcome as x = p (a, s, q, e). Both the owner and management are utility maximizes. The utility function of the manager is formulated on outcome (x) and effort (e), that is. $U_m$ $(x_m,e)$. Management's utility function increases with respect to wealth $(\partial U_m()/\partial x_m > 0)$; it decreases and is strictly concave with respect to effort $(\partial U_m(0)/\partial e < 0$ and $a U_m(0)/\partial x_m)$; management may be either risk neutral $(\partial^2 U_m()/\partial x_m^2)$, or risk-averse $(\partial^2 U_m(0)/\partial x_m^2)$. The owner on the other hand, is concerned only with maximizing his/her residual value less information costs $(x_n)$ – that is, $x_0 = x - x_m - x_n$ where $(x_0)$ is the owner's share from the output-and he is risk neutral $(U_0(x_0) = x_0)$.

It is assumed that the control mechanism is designed in such a manner that signals generated by the system $(n_3)$ are jointly observed by both individuals and, therefore, can be expended as a control element. In this case, the domain of the contract is limited to $x_m = I[\hat{n}_3(s,a,q,e)]$. Given $(I, n_3)$, labor's behavior can be represented by:

$$E(U_m \mid s,a,e,q) = \int U_m[I(\hat{n}_3(a,s,q,e),a]P(s)ds \qquad (6)$$

Let (e| (I,a)) denote management's optimal effort, given a specific contract, (I) and his/her action; then we have:

$$E[U_m \mid e(I), I] = \operatorname{Max} E(U_m \mid e, I) \qquad (7)$$

The owner's utility, on the other hand, depends on (I) and his/her prediction of management's optimal level of effort, $e^*(I)$; thus, his/her expected value model can be expressed as:

$$E(U_0 \mid e^*(I), I) = \iint U_0\{(P(a,s,e^*(I),q) - I(\hat{n}_3(a,s,e^*(I),q) -$$

$$C[(a,s,e(i),q,\hat{n}_3)]\} .P(s)ds.P_e(a|I)da \qquad (8)$$

Where (c) denotes the information cost function; (a), (s), (q), $(e^*)$ and (x) are scalars; and the outcome function has the special multiplicative form, $x = P(a,s,q,e)g(s)\psi F(a,q)$, where g(s) is the stochastic component, $\Psi$ is a positive scalar representing the agent's action, and F(a,q) is homogeneous, increasing, concave, and differentiable with respect to the designated elements.

To determine the type of optimal incentive contract that should be established between the owner and management, Demski and Feltham (1978) first introduced a new kind of contract called a "budget-based contract" with the following characteristics:

1) The agent's compensation is, in part, a function of some observable attribute(s) of the outcome resulting from his/her action.
2) The contract specifies a budgeted (standard) outcome (attribute) level that partitions the set of possible outcomes into favorable and unfavorable subsets.
3) The agent's compensation consists of two functions, one defined over the favorable subset and the other over the unfavorable subset (p. 337).

Having defined the significant characteristics of the budget-based contracts, they considered the implications of informational asymmetry problems. In particular, they investigated situations it which budget-based contracts ought to be established between the owner and management when moral hazard and adverse selection effects

**Table 2.** Characteristics of Pareto-optimal contracts for the moral hazard case.

| A. Costless action or state information | B. Effort and state not observed |
|---|---|
| **Proposition 1** | **Proposition 4** |
| If management's action can be observed without a cost by the risk-neutral owner, then the competitive equilibrium based on the wage contracts will be Pareto optimal. | A necessary condition for a budget based contract to be Pareto-superior to all alternative contracts in the basic moral hazard is that the agent be risk averse. |
| **Proposition 2** | **Proposition 5** |
| If the state can be observed by both the risk- neutral owner and manager without a cost, and the agent has sufficient wealth to provide insurance, then the rental and insurance contracts will be Pareto-optimal. | Given constant stochastic returns for agent and Cobb-Douglas production functions and assuming interior effort solutions, there always exists a dichotomous contract that is Pareto superior to a linear contract. |
| Proposition 3 | **C. Effort or state observed at a cost** |
| | **Proposition 6** |
| A necessary condition for a budget- based contract to be Pareto-superior to all alternative contracts in the basic moral hazard model is that it be costly to observe both labor's effort and the state. | If the principal is risk neutral and the agent is risk averse, then there exists a budget-based contract that is Pareto superior to the linear contracts. |

**Table 3.** Characteristics of Pareto-optimal contracts for the adverse selection case.

| A. Costless skill or state information | B. Skill and state not observed: |
|---|---|
| **Proposition 1** | **Proposition 4** |
| If agent's skill and effort can be observed without a cost by the risk neutral owner, then the competitive equilibrium based on the skill-dependent wage contracts will be Pareto- optimal. | A necessary condition for a budget- based contract to be Pareto-superior to all alternative contracts in the basic adverse selection model is that it be costly to observe both agent's skills and state. |
| **Proposition 2** | **Proposition 5** |
| If the state and capital can be without a cost observed by both the risk neutral owner and manager, and the owner has sufficient wealth to provide insurance, then the competitive equilibrium based on rental and insurance contracts will be Pareto-optimal. | Given certain economic assumptions and with each agent's action, capital, and outcome without a cost observed, there exists a dichotomous budget contract which is Pareto superior to the separating linear contracts. |
| | **C. Effort of State Observed at a Cost** |
| **Proposition 3** | **Proposition 6** |
| A necessary condition for a budget- based contract to be Pareto-superior to all alternative contracts in the basic adverse selection model is that it be costly to observe both agent's skill and the state. | If the principal is risk neutral and the agent is risk averse, then there exists a budget-based contract that is Pareto superior to the linear contracts. |

are expected. A summary of the results of this study concerning the moral hazard selection efforts is shown in Table 2 (Namazi, 1985: 130).

Table 2 lists the types of Pareto-optimal contracts that can be offered to management in the presence of the moral hazard under each observability assumption. For example, given proposition 1, wage contracts are efficient. Rental and insurance contracts provide a first-best solution under proposition 2. Budget-based contracts with costly investigations Pareto-optimally dominate linear contracts under proposition 6.

Similarly, the ADVERSE selection effect case is shown in Table 3 (Namazi, 1985: 138). Hence, these findings

can help managerial accountants select an optimal contract from among different modes of alternative contracting, and they provide a useful theoretical basis for designing managerial budgets and control systems.

## DISCUSSION AND CONCLUSIONS

The major purpose of this study was to extricate the role of the agency theory in implementing management's control system. By adapting the agency theory paradigms and contractual agreement frameworks, it was demonstrated that agency theory has posited, at least, the

following 6 vital roles in this realm:

1) It provides a delegate and precise quantitative and scientific mathematical model to explain why control is important, and therefore, should be exerted in the organization. In effect, it offers a convincing explanation for implementing this powerful managerial accounting technique. This is the "organizational role" of the agency theory.

2) It examines different "performance measures" that must be encompassed in a control system in order to attain a suitable performance. Thus, it solves one of the significant management accountants' obstacles-i.e. what to choose as a performance measures for controlling operations and/or, rewarding the stakeholders. This is part of the "system evaluation" role of the agency theory.

3) It establishes appropriate standards of performance and how are they must be implemented in an attempt to attain Pareto optimality in the organization. Thus, it provides an efficient resource allocation mechanism for the firms. This function is the "allocation efficiency" role of the agency theory.

4) It unambiguously and mathematically demonstrates the significance of various information, and particularly accounting information, in establishing control systems under various conditions of uncertainty. This is the "informational efficiency" role of the agency theory.

5) It encompasses various constituent factors affecting implementation of the efficient control systems. Consequently, the role of pecuniary and non-pecuniary, and behavioral aspects is particularly revealed. This is the "behavior aspects" of the theory.

6) It leads various stakeholders (managers and owners) on how to select a suitable type of contractual agree-ments in different situations, and provides an optimal control mechanism for each realm. This is the "optimal control monitoring selection" aspect of the agency theory.

Despite of the preceding advantages, the presented agency paradigm was basic, primitive and simple. Recent developments in the domain of the agency theory, (Eldenburg and Krishnan, 2003; Kornish and Levine, 2004), however, it has been extended it to situations in which private information, multiple agents, multiple prin-ciples, and multiple objectives are presumed. In addition, the major assumptions of the agency theory, particularly the self-interest behavior, work-aversion, and shirking attitude of the agent, have been questioned, and have entered to be scrutinized under different cultures (Osterman, 2006; Kren and Tyson, 2009). The emergence of empirical research in the agency area has also created a rethinking attitude concerning the agency theory, and has enhanced the existing knowledge in this area, by providing empirical evidence (Dikolli, 2001; Stevens and Thevaranjan, 2010; Demski et al., 2009). While some of these studies have confirmed the agency premises, others have reported the opposite. Hence, more research in this domain, at the theoretical and empirical realm, is suggested for future control studies and endeavors.

## ACKNOWLEDGEMENTS

## REFERENCES

Alchian AA, Demsetz H (1972). Production, information costs and economic organization. Am. Econ. Rev. pp.777-795.

Anthony RN (1965). Planning and control systems: A framework for analysis. Boston: Harvard Graduate School of Business.

Arrow KJ (1971). Political and economic evaluation of social effects and externalities. in Frontiers of quantitative economics, ed. By M. Intriligator (North-Holland).

Baiman S (1982). Agency research in managerial accounting: a survey. J. Account. Lit. (1):154-213.

Baiman S (1990). Agency research in managerial accounting: A second look. Accounting, Organization and Society. 15(4):341-371.

Demski JS (1980). A simple case of indeterminate financial reporting. Information Economics and Accounting Research, ed. By G. Lobo and M. Maher (University of Michigan: Ann Arbor).

Demski JS, Dye RA (1999). Risk, return, and moral hazard. J. Account. Res. 37(1):27-55.

Demski JS, Feltham G (1976). Cost determination: A conceptual approach (Ames: Iowa State University Press).

Demski JS, Fellingham JC, Lin HH (2009). Tension relevance. J. Manage. Account. Res. (21):241-248.

Diamond DW, Verrecchia RE (1982). Optimal managerial contracts and equilibrium security prices. J. Financ. pp.275-287.

Dikolli SS (2001). Agent employment horizons and contracting demand for forward-looking performance measures. J. Account. Res. 39(3):481-494.

Eldenburg L, Krishnan R (2003), Public versus private governance: A study of incentives and operational performance. J. Account. Econ. pp.377-404.

Eisenhardt KM (1989). Agency theory: An assessment and review. Acad. Manage. Rev. 14(1):57-74.

Gjesdal F (1981). Accounting for stewardship. J. Account. Res. Spring:208-231.

Harris M, Raviv A (1976). Optimal incentive contracts with imperfect information. Working paper, Carnegie-Mellon University.

Harris M, Raviv A (1978). Some results on incentive contracts with application to education and employment, health insurance, and law enforcement. Am. Econ. Rev. pp.20-30.

Holmstrom BR (1979). Moral hazard and observability. Bell J. Econ. Spring:74-91.

Horngren CT, Datar SM, Rajan M (2012). Cost accounting: A managerial emphasis", 14th edition, Prentice Hall.

Jensen MC, Meckling WH (1976). Theory of the firm: managerial behavior, agency costs and ownership structure. J. Financ. Econ. pp.305-360.

Kaplan R (1984). The evolution of management accounting. Account. Rev. pp.390-418.

Kaplan R, Atkinson AA (2012). Advanced management accounting. Prentice Hall.

Kornish LJ, Levine CB (2004). Discipline with common agency: the case of audit and non-audit services. Account. Rev. 79(1):173-200.

Kren L, Tyson T (2009). Trade-offs in objective and subjective performance evaluation: a case study examining the validity of the agency theory prediction. Manage. Account. Quart. pp.12-23.

Lan L, Heracleous L (2010). Rethinking agency theory: The view from law. Acad. Manage. Rev. 35(2):294-314.

Namazi M (1983). Accounting information and optimal incentive

contracts in a multi-objective setting. Ph.D. dissertation, University of Nebraska, Lincoln. USA.

Namazi M (1985). Theoretical developments of principal-agent employment contract in accounting: the state of the art. J. Account. Lit. 4:113-163.

Osterman P (2006). Overcoming oligarchy: culture and agency in social movement organizations. Adm. Sci. Quart. 51(4):622-649.

Outley D (1995). Management control, organizational design and accounting information systems. Published in Issues in management accounting, 2nd Edition, Edited by Ashton D, Hopper T, and Scapens RW, Prentice-Hall Europe.

Pacharn P (2008). Accounting choice and optimal incentive contracts: a role of financial reporting in management performance evaluation. Adv. Manage. Account. 171:289-316.

Penno M (1984). Asymmetry of Pre-decision information and management accounting. J. Account. Res. Spring:177-191.

Radner R (1981). Monitoring cooperative agreements in a repeated principal-agent relationship. Econometrica pp.1127-1148.

Ross SA (1973). The economic theory of agency: the principal's problem. Am. Econ. Rev. pp.134-139.

Shavell S (1979). Risk-sharing and incentives in the principal-agent relationship. Bell J. Econ. Spring:55-73.

Spence M, Zeckhauser R (1971). Insurance, information, and individual action. Am. Econ. Rev. pp.380-391.

Stevens DE, Thevaranjan A (2010). A moral solution to the moral hazard problem. Account. Organ. Soc. 35(1):125-139.

Wilson RB (1968). The theory of syndicates. Econometrica pp.119-132.

Zimmerman JL (2010). Accounting for decision making & control. 7th edition, McGraw-Hill; Irwin.

# Accepted Manuscript

Please cite this article as: Liu, Yu, Miletkov, Mihail K., Wei, Zuobao, Yang, Tina, Board independence and firm performance in China, *Journal of Corporate Finance* (2014), doi: 10.1016/j.jcorpfin.2014.12.004

# Board independence and firm performance in China

Yu Liu
College of Business Administration
University of Texas at El Paso, El Paso, TX


Mihail K. Miletkov*
Paul College of Business and Economics
University of New Hampshire, Durham, NH


Zuobao Wei
College of Business Administration
University of Texas at El Paso, El Paso, TX


Tina Yang
School of Business
Villanova University, Villanova, PA

(This version: November 25, 2014)


## Abstract

We provide the first comprehensive and robust evidence on the relationship between board independence and firm performance in China. We find that independent directors have an overall positive effect on firm operating performance in China. Our findings are robust to a battery of tests, including endogeneity checks using instrumental variables, the dynamic generalized method of moments estimator, and the difference-in-differences method. The positive relationship between board independence and firm performance is stronger in government controlled firms and in firms with lower information acquisition costs. We also document that Chinese independent directors play an important role in constraining insider self-dealing and improving investment efficiency.

* Corresponding author: Mihail K. Miletkov, Paul College of Business and Economics, University of New Hampshire, Durham, NH, USA. Tel: 1-603-862-3331; Email: mihail.miletkov@unh.edu

## 1. Introduction

Corporate board structure and its effect on firm decisions and performance is one of the most researched areas in contemporary corporate finance. However, despite the voluminous research in this area, which spans more than two decades and primarily studies U.S. firms, there is no clear evidence of a robust relationship between board composition and firm performance (Hermalin and Weisbach, 2003). This lack of a causal relationship does not necessarily imply that board structure is irrelevant. Rather it is consistent with the view that internal governance mechanisms, such as board structure, are endogenously determined and represent efficient responses to firms' contracting and operating environments (Linck, Netter, and Yang, 2008; Wintoki, Link, and Netter, 2012, among others).

In stark contrast to the U.S. based evidence, studies that use non-U.S. data have consistently documented a positive relation between board independence and firm performance, indicating a possible substitution effect between internal and external governance mechanisms.[1] Compared to the U.S., countries in these studies typically have less-developed legal and extra-legal institutions to provide protection for investor rights. In these countries, internal governance mechanisms, such as board structure, become more consequential (Klapper and Love 2004; Ferreira and Matos, 2008; McCahery, Sautner, and Starks, 2010).

In this paper we study the impact of board independence on firm performance in the listed firms of the most populous country and second largest economy in the world—China. It is unclear whether the positive relationship between board independence and firm performance documented in

---

[1] Multinational studies include Dahya, Dimitrov, and McConnell (2008), Aggarwal, Erel, Stulz, and Williamson (2009), and Bruno and Claessens (2010). Single country studies include Yeh and Woidtke (2005) for Taiwan, Black and Khanna (2007) for India, Dahya and McConnell (2007) for the U.K., and Choi, Park, Yoo (2007) and Black and Kim (2012) for Korea.

1

several non-U.S. countries can be generalized to China due to the unique characteristics of its listed sector. The most distinctive feature is that the majority of Chinese listed firms are former state-owned enterprises (SOEs) and that the government is still the largest shareholder in many of those companies. On the other hand, China is similar to many other emerging-market countries in that it has a concentrated corporate ownership structure and a weak institutional environment where investor protection is poor and insider self-dealing is rampant (Jiang, Lee, and Yue, 2010).

We use a data set covering almost all publicly traded firms on the Shanghai and Shenzhen stock exchanges from 1999 to 2012 to examine the relationship between board independence and firm performance. Our results show that board independence improves firm operating performance in China. We further document that the positive effect of board independence on performance is more pronounced in government-controlled firms and in firms whose independent directors face lower costs of acquiring firm-specific information.

Our study is largely motivated by the current state of the literature on board independence and firm performance in China. In a recent survey paper, Wang (2014) summarizes the evidence from 30 empirical studies investigating the relationship between board independence and firm performance in China. Wang (2014) documents mixed results, with five (four) papers reporting a significantly positive (negative) relationship while the remaining twenty-one reporting an insignificant relationship. Wang (2014, page 5) summarizes the current state of the literature: "[E]mpirical research on the association between independent directors and firm performance in China seems to be abundant in scope but not plentiful in depth." In our study, we address this deficiency in the literature by providing the first in depth and comprehensive analysis on the relation between board independence and firm performance in China's listed sector. Specifically, we employ several identification strategies,

2

including firm fixed-effects, difference-in-differences (DID), two-stage least squares with instrumental variables (IV-2SLS), and dynamic generalized method of moments (GMM) regressions, to explicitly account for potential endogeneity concerns.

More importantly, most of the studies cited in Wang (2014) simply examine the relationship between board independence and firm performance in China's listed firms. We advance the literature by going beyond that—we investigate the "how," i.e., the potential channels through which independent directors exert their influence on firm operations. Prior literature establishes that insider self-dealing and government intervention are major impediments to efficient operation and investment in Chinese listed firms (see, e.g., Jiang, Lee, and Yue, 2010; Chen, Sun, Tang, and Wu, 2011; Wang and Xu, 2011). We extend this literature by examining the effect of board independence on the tunneling of firm resources through intercorporate loans (a main tool for insider self-dealing in China) and on the firm's investment behavior and outcomes. We find that board independence reduces tunneling through intercorporate loans and improves investment efficiency, especially in government-controlled firms.

According to China's Code of Corporate Governance of Listed Companies, shareholders with more than 1% of total outstanding shares can nominate independent directors.[2] In fact, the Shanghai Stock Exchange (2004) reports that 70% of the independent directors are nominated by top shareholders of the firms. We find that in government-controlled firms board independence reduces insider self-dealing, improves investment efficiency, and enhances firm performance. Therefore, our findings support the notions that independent directors are effective monitors and that the government

---

[2] The Code was jointly issued by the China Securities Regulatory Commission (CSRC) and the Chinese National Economic and Trade Commission in January 2002. Please see Appendix 1 for a comparison of the regulations and characteristics of independent directors between China and the U.S.

3

appoints stronger boards to prevent insider self-dealing and to credibly signal a commitment not to interfere in company affairs.

Our findings have important policy implications. Independent directors became a prevalent feature in China's corporate governance landscape only following the issuance of *The Guidelines for Introducing Independent Directors to the Board of Directors of Listed Companies*. These guidelines, which were introduced in 2001 by the Chinese Securities Regulatory Commission (CSRC), mandated that by June 30, 2003 the boards of all listed firms should be comprised of at least one-third independent directors. Using this regulatory mandate as an exogenous shock, our DID analysis shows that the mandatory board independence provision is effective in improving firm performance in China's listed firms. Further analyses show that firms, which voluntarily appointed independent directors prior to the mandate, exhibit higher performance than those that did not. In addition, firms that appointed more independent directors than the required minimum following the mandate exhibit significantly better performance than those who choose not to go over the required minimum. Taken together, our results are consistent with the prevailing evidence from other emerging markets (e.g., Dahya, Dimitrov, and McConnell, 2008; Black and Kim, 2012) that good internal governance practices, such as independent boards, are effective at curbing agency problems, including those associated with dominant shareholders, thereby leading to better performance. In this light, the evidence from our study offers additional empirical support for the global push by regulators and governance activists for good corporate governance practices (Aggarwal, Erel, Ferreira, and Matos, 2011).[3]

The remainder of the paper proceeds as follows. In Section 2 we review the related literature and summarize our research design. Section 3 presents our sample and variable descriptions. In

---

[3] Also see Section 2 for more details on the regulatory push in different countries for more independent boards.

4

Section 4 we discuss our main empirical results on the relationship between board independence and firm performance. In Section 5 we go beyond the independence-performance relation and investigate the potential channels through which independent directors influence firm performance in China. Section 6 presents results from additional robustness checks, including those on the effect of board independence beyond the 2001 regulatory mandate, and Section 7 provides concluding remarks.

## 2. Related literature and research design

Hermalin and Weisbach (2003) survey the literature on the relationship between board composition and firm performance in U.S. firms, and draw the conclusion that board composition does not impact firm performance. They further argue that corporate board structure is endogenously determined and those studies found a significant relationship between board composition and firm performance because they failed to adequately control for the endogenous relationship.

In this context, it may seem surprising that lawmakers and stock exchanges around the world have embraced board independence as an essential element of "good" corporate governance and have adopted legislation or codes prescribing higher representation of outsiders on the boards of publicly traded companies.[4] The desirability of mandating more independent boards, however, accords well with the empirical evidence from most studies on the efficacy of independent boards outside the U.S. For example, in their study of 22 non-U.S. countries, Dahya, Dimitrov, and McConnell (2008) document that board independence is significantly and positively related to firm performance,

---

[4] The following is a list of countries that have adopted a minimum standard for outsider board representation as reported in Dahya and McConnell (2007) (page 540): Australia, Belgium, Brazil, China, Cyprus, Czech Republic, Denmark, France, Greece, Iceland, Indonesia, India, Japan, Kenya, Malaysia, Mexico, New Zealand, Poland, Portugal, Russia, Singapore, South Africa, South Korea, Sweden, Switzerland and Thailand.

5

especially in countries with lower levels of investor protections. These results are further corroborated by the findings in Aggarwal, Erel, Stulz, and Williamson (2009) and Bruno and Claessens (2010).

Several single country studies also document a positive effect of board independence on firm performance. Black and Khanna (2007), Dahya and McConnell (2007), and Black and Kim (2012) all examine country-specific regulatory shocks and conclude that increasing board independence significantly improved firm performance in India, the U.K, and Korea, respectively.

Researchers have provided several explanations for the contrasting evidence between U.S. and non-U.S. firms concerning the effect of board independence on firm performance. First, there may be a substitution effect between internal and external governance mechanisms (Klapper and Love, 2004; Ferreira and Matos, 2008; McCahery, Sautner, and Starks, 2010) such that monitoring by outside directors is much more important in those non-U.S. countries where legal and extra-legal institutions offer substantially weaker investor protections.[5] Additionally, in the U.S. most corporate boards have long been dominated by outside directors. Cicero, Wintoki, and Yang (2013) report that between 1991 and 2003 the level of board independence in the average U.S. publicly traded firm ranged from 63% to 71%. This lack of significant variation in the degree of board independence across firms and through time may also preclude the identification of a statistically significant relationship between board independence and firm performance in the U.S. data alone.

In this paper we examine the effect of board independence on firm performance in Chinese listed firms. While prior studies have examined this research question, they fail to identify a robust

---

[5] The extra-legal institutions include, but are not limited to, the market for corporate control, culture and norms, product market competition, monitoring by financial market participants, and the financial press (Dyck and Zingales, 2004).

relationship between board independence and firm performance in China (Wang, 2014). There are three possible explanations for this lack of a robust relationship in the extant literature. First, Chinese corporate boards may be at their optimal construction, and thus, no relationship could be observed in the aggregate. However, it is highly unlikely that the average firm in China's listed sector has already achieved the optimal board structure, as evidenced by the existence of rampant agency problems in China's listed firms (Sun and Tong, 2003; Allen, Qian, and Qian, 2005; Liu, Wei and Xie, 2014). The second possible explanation is that independent directors in China are perfunctory. They are ineffective at monitoring and advising the management and therefore, no robust relation can be observed in the aggregate. If this is the case, the validity of the aforementioned 2001 government mandate for more independent directors would be called into question. This leads to a third, and more likely explanation: studies in the extant literature may have failed to thoroughly account for the endogenous relation between board independence and firm performance arising from potential simultaneity bias, unobserved heterogeneity, and/or reverse causality.[6]

To address the aforementioned econometric challenges that plague most of the extant literature, we construct a comprehensive data set covering almost all publicly listed Chinese firms during the 14-year period from 1999 to 2012. We then employ several identification strategies specifically designed to address potential endogeneity in the relationship between board independence and firm performance. Given the mixed results from prior empirical studies and the unique characteristics of the Chinese corporate landscape, we make no prediction regarding the overall effect of board independence on firm performance.

---

[6] Similar arguments have been put forth to explain why empirical researchers have failed to detect a robust relationship between board independence and firm performance for U.S. firms (see, e.g., Hermalin and Weisbach, 2003; Dahya and McConnell, 2007; Adams, Hermalin, and Weisbach, 2010).

7

Throughout the paper we use the fraction of independent directors on corporate boards as the measure of board independence. However, in additional analysis we also investigate whether the actual *number* of independent directors affects firm performance. This additional analysis is important because while most of the literature focuses on the fraction of independent directors, many of the recommendations or mandates issued by governments and/or stock exchanges around the world prescribe a minimum number of outside directors. The first and most widely recognized example is the Code of Best Practice issued by the Cadbury Committee in 1992, which recommends that U.K. listed firms should have at least three outside directors. Dahya and McConnell (2005, 2007) document that since 1992, at least 15 other countries have adopted either voluntary or mandatory standards for a minimum number of outside directors on corporate boards. In most cases, the prescribed minimum number of outside directors is either two or three. We view this analysis as empirical as there is very little theoretical guidance in the literature as to what an "optimal" number of independent directors might be.[7]

As mentioned earlier, the Chinese stock market is unique in that the government is the largest shareholder in most of the listed firms, and the effect of government ownership on the relationship between board independence and firm performance is ambiguous. On the one hand, prior literature suggests that the Chinese government is not a passive shareholder; rather, it actively intervenes in company management and often compels firms to pursue social and political goals in lieu of shareholder wealth maximization (Fan, Wong, and Zhang, 2007; Chen, Sun, Tang, and Wu, 2011). In this context, independent directors may lack the ability and/or incentives to actively monitor and

---

[7] However, in the context of board gender and race diversity, several researchers have documented the existence of a critical mass (Broome, Conley and Krawiec, 2011; Liu, Wei, and Xie, 2014) (see more discussions in Section 4.2).

8

discipline the management. On the other hand, the government may use the appointment of independent directors to signal a commitment not to interfere in company affairs, similar to the way controlling shareholders in countries with low investor protections appoint strong (i.e., more independent) boards in order to convince outside investors that they will refrain from diverting corporate resources (Dahya, Dimitrov, and McConnell, 2008). Furthermore, independent board members could be especially valuable as monitors in government-controlled firms that are subject to severe agency problems due to the ultimate separation between ownership (the nation's citizens) and control (managers/bureaucrats).[8] Indeed, there is some anecdotal evidence that independent directors in government-controlled firms are pulled in different directions. For example, Xiangbin Yin, a Chinese independent director, described his experiences as follows: "[T]he State, the largest shareholder, wants him to be a 'KGB' in the company to ensure the integrity of the managers; while minority shareholders expect him to be a 'white knight' to fight against the exploitation from the controlling shareholder and from managers (Shen and Jia, 2005, page 233)".

Prior research on U.S. firms further shows that independent directors are more valuable in firms with lower information acquisition and monitoring costs (Linck, Netter, and Yang, 2008; Duchin, Matsusaka, and Ozbas, 2010). We extend the literature by examining whether in China's listed firms the independence-performance relation is also a function of firm-specific information acquisition and monitoring costs.

In addition to providing robust evidence on the effect of board independence on firm performance in China's listed firms, our paper makes another important contribution by examining the

---

[8] See Megginson and Netter (2001) for a survey of the literature on the effects of government ownership on firm behavior and performance.

9

potential channels through which independent directors can exert their influence on firm operations and performance. If independent directors are effective at monitoring managers and preventing malfeasance, we should expect to find a negative relationship between board independence and insider self-dealing. To test this hypothesis, we follow Jiang, Lee, and Yue (2010) and use the magnitude of intercorporate loans as a proxy for insider self-dealing. We further split our sample into two time periods, 1999-2005 and 2006-2012, to account for the fact that the practice of tunneling through intercorporate loans was largely eliminated by 2006 following a series of government directives and actions (Jiang, Lee, and Yue, 2010). We expect that any effect of board independence on this form of insider self-dealing should be observed in the pre-2006 period.

We also examine the effect of board independence on firms' investment behavior. Prior literature suggests that agency problems are some of the main factors explaining why firm investment may deviate from its optimal level (Jiang, Kim, and Pang, 2011). Therefore, if monitoring by independent directors reduces agency conflicts, we should expect the sensitivity of investment expenditures to investment opportunities to be higher in firms with more independent boards. Furthermore, as documented by Chen, Sun, Tang, and Wu (2011), government intervention in firms with majority state ownership also leads to investment inefficiency, as those firms are being pressured to pursue politically expedient projects, rather than value-maximizing ones.[9] If the appointment of independent directors reduces the likelihood of government interference in company affairs, then we should expect to observe a positive relationship between board independence and investment efficiency in government-controlled firms.

---

[9] Politically expedient projects include projects that maximize employment, promote regional development, foster social stability, etc. (Chen, Sun, Tang, and Wu, 2011).

10

### 3. Sample and variable descriptions

*3.1. Sample*

We start the sample collection process with all the listed firms on the Shanghai and Shenzhen Stock Exchanges for the period of 1999 to 2012. We obtain data on financial statements, stock prices, board composition, and ownership structure from the Chinese Securities Market and Accounting Research (CSMAR) Database, a leading provider of data on Chinese stock markets and listed companies. We exclude from the sample financial companies and utilities, as well as firms with negative or zero net assets. To mitigate extreme outliers, we truncate firm performance measures by 1% at both tails. The final sample with the requisite financial, board, and ownership data consists of 16,999 firm years or 2,057 unique firms.

*3.2. Measures of firm performance*

To measure firm performance, we choose accounting measures, specifically, return on assets (ROA) and return on equity (ROE), over market-based measures such as Tobin's Q or stock returns because Chinese stock ownership underwent significant reform during our sample period, rendering accounting measures much more responsive to the underlying firm economic performance than market-based measures. Specifically, Chinese stock markets were re-opened in the early 1990s with the establishment of the Shanghai Stock Exchange in 1990 and the Shenzhen Stock Exchange in 1991. To facilitate a gradual and smooth transition from a planned economy to a market economy, newly privatized, listed companies were required to have multiple classes of shares: (1) "A-shares" denominated in Chinese currency, RMB, and designated for domestic investors, (2) "B-shares" denominated in U.S. or HK dollars, and reserved for foreign investors or domestic investors with foreign currencies, and (3) "H-shares" for those that are cross-listed in Hong Kong Stock Exchange. A-shares were further divided into tradable and non-tradable shares. Tradable A-shares were the only class of shares that could be traded among domestic investors and they have turned out to be highly speculative subject to extreme turnover rates (Mei, Scheinkman, and Xiong, 2009; Liu, Wei, and Xie,

11

2014). Non-tradable shares, which entitled the holders to exactly the same rights as holders of tradable

A-shares except for public trading, were typically held by the government or other institutions. Before

2005, two-thirds of the Chinese stock market was comprised of non-tradable shares. In 2005, the

CSRC required all listed firms to transform their non-tradable shares into tradable shares. The reform

was largely completed by 2008, subject to certain lockup provisions and trading limits.

We compute ROA (ROE) as operating income before extraordinary items divided by total

assets (total equity). As Panel A of Table 1 shows, the mean values of ROA and ROE are 3.6% and

5.7%, respectively, over the period of 1999 to 2012, which is in line with the existing literature. For

example, Jiang, Lee, and Yue (2010), who also calculated ROA as pre-extraordinary operating income

over total assets, but winsorized the variable by 1% at both tails, report a mean ROA of 2.8% with a

standard deviation of 6.3%.

### 3.3. Measures of board composition

We use the percent of independent directors (*%_Ind*) to measure board composition. When

classifying independent directors, we follow the guidelines of the CSRC. Specifically, to qualify as an

independent director, the director cannot be: (1) a current or former employee of the company or its

subsidiaries, (2) a(n) (in)direct owner of more than 1% of the outstanding shares of the company, (3)

the legal person (or direct family members of the legal person) of the top 10 owners of the company,

(4) an employee of an institution that directly or indirectly owns more than 5% of the outstanding

shares of the company, (5) an employee (or direct family member of an employee) of the top five

owners of the company, or (6) a provider of financial, legal, or consulting services to the company and

its subsidiaries. Lastly, independent directors need to meet the independence requirements specified in

the charter and bylaws of their respective companies.

We also use the actual number of independent directors to investigate whether it has an effect

on overall firm performance. Specifically, the indicator variables *Ind_d1*, *Ind_d2*, *Ind_d3*, and *Ind_d4*

equal one when the board has one, two, three, or four independent directors, respectively, and zero

otherwise. The indicator variable *Ind_d5* equals one if the board has five or more independent

12

directors and zero otherwise. Table 1 Panel B reports the summary statistics for our board independence measures.

Figure 1 depicts the time trend of board independence for our sample period from 1999 to 2012, which can be better understood given the context of the development of corporate governance reform in China. China's corporate governance laws and practices underwent watershed change in 2001. In that year, the scandal involving Ying Guang Xia, dubbed as "China's Enron," came into light. That year also marked the start of a series of reforms to improve corporate governance of Chinese listed companies. From that time on, corporate governance moved from being a virtually unknown concept in China to the centerpiece of Chinese economic reform. In 2001, China joined the World Trade Organization and adopted *the Principles of Corporate Governance* published by the Organization for Economic Co-operation and Development (OECD). On August 16, 2001, the CSRC issued *The Guidelines for Introducing Independent Directors to the Board of Directors of Listed Companies* (hereafter the 2001 Guidelines), requiring that by June 30, 2002, each listed firm in China shall have at least two independent directors and that by June 30, 2003, at least one-third of the board shall be comprised of independent directors. "Independent director" is defined as independent of management and the relatives of the management, the controlling shareholder, and the persons providing financial, legal or consulting services to the company. In January 2002, the CSRC and the Chinese National Economic and Trade Commission jointly issued *The Code of Corporate Governance of Listed Companies* (the 2002 Code) to speed up the governance reform.

Consistent with the series of regulatory pushes for higher board independence, we observe a dramatic increase in board independence starting in 2001. For example, the percent of firms with three or more independent directors increased from 1% in 1999 to 85% in 2003 and then to 97% in 2012. Notably, the trend of board independence increased at a much slower pace post 2003. For example, the percent of independent directors rose from 1% in 1999 to 33% in 2003, but by only 4% from 2004 to 2012. As we discuss in Section 4.1.2, the implementation of the 2001 Guidelines created an exogenous shock that allows us to empirically identify the effect of board independence on firm performance.

13

[Insert Figure 1 here]

*3.4. Control variables*

The control variables that we include in all performance regressions are grouped into three broad categories: ownership variables, proxies for monitoring costs, and others. Summary statistics for these variables are presented in Table 1 Panel C.

The top owner of a Chinese listed firm can be: 1) the State or a state-owned enterprise (SOE); 2) one or several individuals; and 3) a non-state-owned legal entity such as another public or private enterprise. We control for ownership structure using *Topowner_State*, a dummy variable that takes the value of one if the top owner of a listed firm is the State or a state-owned enterprise (SOE), and *Topowner_Individual*, a dummy variable that takes the value of one if the top owner is one or several individuals. We also control for foreign ownership with *%_Foreign*, the percent of B-shares and H-shares issued by the firm. Following the literature (see, e.g., Fama and Jensen, 1983; Linck, Netter, and Yang, 2008), we use two variables to proxy for monitoring costs: (1) stock price volatility (*Volatility*), which is the annualized standard deviation of weekly stock returns, and (2) the geometric mean of sales growth rate over the past three years (*Sales_Growth*).[10]

We also include other commonly used control variables for board and firm characteristics, i.e., the natural log of the number of directors on the board (*Ln_BoardSize*), an indicator variable for whether the Chief Executive Officer chairs the board (*Duality*), investment expenditure divided by total assets (*Invest_Expenditure*), the natural log of total assets (*Ln_Assets*), total liabilities divided by total assets (*Leverage*), and the natural log of the number of years that a firm has been listed on an exchange (*Ln_FirmAge*).

[Insert Table 1 here]

---

[10] Unfortunately, our data sources do not provide information on some of the other proxies used in the literature for monitoring costs, such as the number of business segments or research and development expenses.

14

In Table 2, we report the correlation matrix of all the variables employed in the subsequent regression analyses. None of the independent variables have correlation coefficients higher than 0.5 in absolute terms, thus alleviating potential concerns about multicollinearity.

[Insert Table 2 here]

## 4. Main results

### 4.1. Do independent directors affect firm performance?

Like most empirical corporate finance research, the analysis of the relationship between board composition and firm performance faces the challenge of endogeneity, which can arise from unobserved heterogeneity, simultaneity, and reverse causality (Adams, Hermalin, and Weisbach, 2010; Wintoki, Linck, and Netter, 2012). In the context of the board-performance relationship, the problem of unobserved heterogeneity arises when one or more latent variables drive the observed relationship between board characteristics and firm performance. For example, high-ability CEOs may be more inclined to choose more independent boards and those CEOs deliver better firm performance. The problem of simultaneity can arise when firms choose certain levels of board independence with the goal of achieving certain levels of firm performance. The problem of reverse causality can arise when firms with a certain level of economic performance are prone to adopting certain levels of board independence.

To robustly test the impact of independent directors on firm performance, we design and implement four empirical tests that use the following estimation methods in a panel regression framework: (1) firm fixed-effects (FE), (2) difference-in-differences (DID), (3) two-stage least squares with instrumental variables (IV-2SLS), and (4) dynamic generalized method of moments (GMM). The firm FE model addresses endogeneity due to unobserved, time-invariant heterogeneities, but does not solve the endogeneity problems due to time-varying heterogeneities, simultaneity, or reverse causality. Therefore, we use DID, IV-2SLS, and dynamic GMM to further address these

15

econometric challenges (Wintoki, Linck, and Netter, 2012; Semadeni, Withers, and Certo, 2013). Table 3 presents the results from all four econometric models applied to our panel data set.

### 4.1.1. Firm fixed-effects estimation

Table 3 columns (1) and (5) report the regression results from estimating the following ordinary-least squares (OLS) model with firm fixed-effects:

$$Performance_{it} = \gamma * \%\_Ind_{it} + \beta * Control_{it} + d_i + d_t + \varepsilon_{it} \qquad \dots \text{Eq. 1}$$

Firm performance (*Performance*) is measured by ROA or ROE. The proportion of independent directors on the board (*%_Ind*) is our key variable of interest. *Control* is a vector of firm-specific variables previously described in Section 3.4 that prior literature has identified as potentially related to firm performance (see, e.g., Wintoki, Linck, and Netter, 2012). $d_i$ and $d_t$ denote firm and year fixed-effects respectively, and $\varepsilon_{it}$ is the error term. We use robust standard errors that are adjusted for firm-level clustering to control for potential serial correlation and heteroscedasticity in the data (Petersen, 2009).

As columns (1) and (5) illustrate, board independence (*%_Ind*) is associated with significantly higher firm performance. An increase of ten percentage points in board independence—equivalent to adding one director to the board given the average board size of 9.6—leads to a statistically significant increase of 29 basis points in ROA or a 138-basis-point increase in ROE. This effect is also economically meaningful as the mean values of ROA and ROE in our sample are 3.6% and 5.7%, respectively. We find some evidence that smaller boards are also associated with better firm performance, which corroborates the findings in prior U.S. studies (e.g., Yermack, 1996). Consistent with the literature (e.g., Qi, Wu, and Zhang, 2000; Guo and Jiang, 2013), we document that state ownership (*Topowner_State*) is significantly and negatively related to firm performance. Most of the other control variables also have the predicted signs.

16

*4.1.2. Difference-in-differences estimation*

DID is a widely used and effective empirical method in the economics and finance literature to address the endogeneity problem (see, e.g., Bertrand and Mullainathan, 2003; Yang and Zhao, 2014). In the context of our study, we use the exogenous shock of a regulatory mandate to analyze whether two groups of firms—one impacted by the mandate (the treatment group) and the other not impacted (the control group)—exhibit different performance trends surrounding the regulatory shock. This design mitigates the endogeneity problem because the change in board composition is triggered by regulation. In other words, the change is not endogenously driven by firm-specific characteristics. As discussed in Section 3.3, in 2001, the CSRC mandated that all listed firms must have at least one-third of their boards comprised of independent directors by June 30, 2003. Figure 1 in Section 3.3 illustrates that the 2001 mandate had a significant impact on board composition of Chinese listed firms. Therefore, the 2001 rule change constitutes the desired natural experiment that we will exploit in the following DID framework to probe the relationship between board independence and firm performance:

$$Performance_{it} = \gamma * Treated_i * Post\_Legislation_t + \beta * Control_{it} + d_i + d_t + \varepsilon_{it} \qquad \dots \text{Eq. 2}$$

*Treated* is an indicator variable that equals one if a firm has less than one-third of independent directors on its board as of 2001 (i.e., *%_Ind<33%*) and zero otherwise. *Post_Legislation* is an indicator variable that equals one if the sample year is 2003 or later and zero otherwise. Simultaneous inclusion of $Treated_i * Post\_Legislation_t$ and firm fixed-effects ($d_i$) removes any biases from comparison between the treatment and the control groups that may result from permanent differences between the two groups. Simultaneous inclusion of $Treated_i * Post\_Legislation_t$ and year fixed-effects ($d_t$) removes any biases from comparison over time in the treatment group that may result from trends. To conduct this test, we require firms to be in the sample before the exogenous shock (i.e., by the end of 2001). This requirement decreases the number of observations from the full sample of 16,999 to 12,604. Columns (2) and (6) in Table 3 report the estimation results from the DID method. The

17

estimated coefficient of *Treated_i*Post_Legislation_t* is significantly and positively related to both ROA and ROE, corroborating our earlier finding that independent boards improve firm performance. As columns (2) and (6) show, ROA and ROE improve 1.4% and 4.7%, respectively, as a result of the mandate for more independent directors.

### 4.1.3. IV-2SLS estimation

IV-2SLS is another standard methodology used to address the endogeneity problem in empirical corporate finance research. A valid instrument must meet two criteria: a strong correlation with the instrumented regressors and orthogonality with the error term. Following the literature, we instrument *%_Ind* of a specific firm in a given year, by using the mean value of the percent of independent directors of other firms in the same industry (one-digit SIC codes) in the same year (*%_IndustryInd*), and the mean value of the percent of independent directors of other firms headquartered in the same province/municipality in the same year (*%_ProvinceInd*). The rationale behind *%_IndustryInd* is that firms' governance arrangements (board composition in our case) likely correlate with their industry peers' due to similar business mix and investment opportunities, but such industry average is unlikely to directly influence individual firm performance (Yang and Zhao, 2014). We also use *%_ProvinceInd* to capture the local supply of directors, which can serve as a valid instrument for board composition for two reasons. First, Knyazeva, Knyazeva, and Masulis (2013) find that the local supply of directors is an important determinant of board independence. Second, the province/municipality in which a firm is headquartered can be viewed as predetermined, because headquarters decisions are made early in the life of a firm and headquarters changes occur very rarely. The suitability of our instruments is confirmed by various identification tests reported in Table 3. In unreported first stage regressions, our instruments significantly correlate with *%_Ind* at a 1%

18

significance level. In the second stage regressions, the Hansen's over-identification test fails to reject

the hypothesis that our instruments are exogenous. As columns (3) and (7) of Table 3 illustrate, the

IV-2SLS method also supports our earlier findings that board independence enhances firm

performance in Chinese listed firms.

### 4.1.4. Dynamic GMM estimation

In an effort to further confirm the robustness of the independence-performance relationship,

we also use the dynamic GMM method to address endogeneity concerns due to unobserved

heterogeneity, simultaneity, and reverse causality. Compared to the DID and IV-2SLS methods,

dynamic GMM has the advantages of: (1) tackling the endogeneity problem based on internal

instruments instead of relying on natural experiments or external instruments, which may not be

readily available, and (2) explicitly modeling the dynamic nature of the independence-performance

relationship by including past firm performance as one of the regressors. Following Wintoki, Linck,

and Netter (2012), we estimate the following regression model:

$$Performance_{it} = \gamma * \%\_Ind_{it} + \alpha * Performance_{i,t-2} + \beta * Control_{it} + d_i + d_t + \varepsilon_{it} \quad \dots \text{Eq. 3}$$

We assume that *%_Ind* and all of the control variables, except for *Ln_FirmAge* and the year

dummies ($d_t$), are potentially endogenous. Following Wintoki, Linck, and Netter (2012), we sample at

two-year intervals instead of every year (i.e., 1999, 2001, 2003…) to minimize serial correlation in

the first-differenced residuals. As a result, two-year lagged ROA or ROE (*Performance_{i,t-2}*) is included

in Equation (3). We also include as IVs the sixth and eighth lags of the dependent and endogenous

variables plus all of the available lags of the exogenous variables. As in previous models, standard

errors are robust to heteroscedasticity and firm-level clustering. To ensure that the dynamic GMM

method is correctly specified, we examine the exogeneity of IVs and the autocorrelation conditions of

the transient errors. As reported in Table 3, the Hansen's over-identification test fails to reject the

19

hypothesis that the selected IVs are exogenous. In addition, residuals are significantly correlated in the first differences (AR(1)), but are uncorrelated in the second differences (AR(2)), suggesting that the assumptions of the dynamic GMM model hold. As columns (4) and (8) in Table 3 demonstrate, board independence has a positive and significant impact on ROA, and a positive but insignificant impact on ROE.

[Insert Table 3 here]

In unreported robustness checks, we redefine ROA (ROE) as net income divided by assets (equity) and re-run all the regressions in Table 3. Our results hold using these alternative measures of firm performance. Many researchers have argued that board variables and other firm characteristics need time to affect firm performance. Therefore, we re-run all the regressions in Table 3 with one year forwarded ROA or ROE (two year forwarded ROA or ROE when using the dynamic GMM method) as the dependent variable. All of our results are robust to these alternative specifications.

## 4.2. Is there a critical mass?

As discussed in Section 2, while most of the extant literature employs the fraction of independent directors as a measure of board independence, many of the recommendations or mandates issued by governments and stock exchanges prescribe a minimum number of outside directors. The 2001 Guidelines mandate that all listed firms in China should have at least two independent directors by June 30, 2002, and a minimum of one-third by June 30, 2003. This suggests that both the number and the fraction of independent directors matter to regulators and potentially to investors as well. In the context of the debate regarding gender and race diversity in corporate boardrooms, there is a belief that unless a "critical mass" is reached, female or ethnic minority directors will not have any meaningful impact on corporate decisions or outcomes (Broome, Conley and Krawiec, 2011). Liu, Wei, and Xie (2014) study the effect of female directors on firm

20

669

performance in Chinese listed firms, and find that when firms have three or more female directors, the impact of board gender diversity on firm performance becomes significant. We, therefore, test whether there is a critical mass of independent directors, and if so, what number constitutes such a critical mass.

In Table 4, we re-estimate Equation (1) by replacing *%_Ind* with *Ind_d1*, *Ind_d2*, *Ind_d3*, *Ind_d4*, and *Ind_d5*, which are binary variables that indicate the presence of one, two, three, four, and five or more independent directors on the board, respectively. Column (1) reports the effects of the independent director dummies on *ROA*. Compared to boards with no independent directors, *Ind_d1* has no significant impact on *ROA*, suggesting that a solo independent director may be a mere token. This is consistent with some of the evidence from psychology literature that a lone individual is unlikely to express a view contrary to an otherwise unanimous group (Asch, 1951). *Ind_d2*, *Ind_d3*, *Ind_d4* and *Ind_d5* are associated with 0.7%, 1.6%, 1.6%, and 1.6% higher *ROAs* than boards with no independent director, respectively. These findings suggest that three independent directors may constitute a critical mass that positively and significantly impacts *ROA*.

Column (2) reports the effects of the independent director dummies on *ROE*. In line with the results in column (1), *Ind_d1* has no significant impact on *ROE*. *Ind_d2*, *Ind_d3*, *Ind_d4*, and *Ind_d5* are associated with 3.6%, 5.4%, 6.5%, and 6.8% higher *ROEs* than boards with no independent director, respectively. These results suggest that ROE improves as the number of independent directors increases.[11]

[Insert Table 4 here]

In our sample, the average board has nine directors, while the median number of independent directors is three, or one-third of the average board. This suggests that most Chinese listed firms appoint just enough independent directors to meet the quota mandated by the CSRC 2001 Guidelines.

---

[11] In an untabulated analysis, we find that the estimated coefficient of *Ind_d2* is significantly different from that of *Ind_d5* in both regressions of *ROA* ($p$-value=0.01) and *ROE* ($p$-value=0.01). The estimated coefficient of *Ind_d3* is not significantly different from that of *Ind_d5* in the regression of *ROA*, but is different from that of *Ind_d5* in the regression of *ROE* with marginal significance ($p$-value=0.10).

21

The results in Table 4 suggest that some firms may be able to further improve their ROE by appointing additional independent directors to their boards.

## 5. When and how independent directors affect firm performance?

In Section 4 we find that independent directors have an overall positive impact on the operating performance of China's listed firms. In this section, we provide four additional analyses examining when and how independent directors can affect firm performance. In Sections 5.1 and 5.2, we test the independence-performance relationship conditioned on the degree of information asymmetry (monitoring costs) and the level of state ownership, respectively. In Sections 5.3 and 5.4, we investigate two key channels through which independent directors can impact firm performance in China, by constraining insider "self-dealing" and by improving investment efficiency, respectively.

### 5.1. Monitoring costs, independent directors, and firm performance

Compared to inside directors, outside directors face at least two disadvantages that can impact the effectiveness of their advising and monitoring functions. The first is the information-asymmetry disadvantage in that inside directors have better information concerning all aspects of the operation and management of the firm. The second is an expertise disadvantage in that inside directors are more likely to possess firm-specific expertise, whereas outside directors are more likely to possess more generic knowledge. Therefore, to be effective in carrying out their advising and monitoring activities, independent directors must incur significant costs related to the acquisition of firm-specific information and expertise. Consistent with this argument, Maug (1997) shows that in high information asymmetry environments, it is suboptimal to increase monitoring by independent directors. Furthermore, Raheja (2005) and Adams and Ferreira (2007) theoretically model board structure and postulate that the number of outside directors declines as the cost of monitoring increases. More

22

671

recently, using a sample of U.S. firms, Duchin, Matsusaka, and Ozbas (2010) directly test and find support for the hypothesis that the benefits of having more independent boards are lower in firms with higher information acquisition costs. We extend this literature by examining whether the moderating effect of information asymmetry on board effectiveness is also present in Chinese listed firms.

Following the literature (Fama and Jensen, 1983; Linck, Netter, and Yang, 2008), we employ two commonly used proxies for information asymmetry/monitoring costs, namely the annualized standard deviation of stock returns (*Volatility*) and the sales growth rate (*Sales_Growth*). We add the interaction terms *%_Ind\*Volatility* and *%_Ind\*Sales_Growth* to our baseline regression, Eq. (1). As Table 5 reports, both interaction terms are negatively and significantly related to firm performance, consistent with the existing literature which finds that as monitoring costs increase, the benefits associated with higher board independence decrease.

[Insert Table 5 here]

*5.2. Ownership structure, independent directors, and firm performance*

As discussed in Section 2, the Chinese stock market is unique in that the government is the largest shareholder in most of the publicly traded companies. We explicitly test the effect of government ownership on the relationship between board independence and firm performance by dividing our sample into government-controlled versus non-government-controlled firms. Table 6 illustrates that the positive relationship between board independence and firm performance, which we document in Table 3, is largely driven by government-controlled firms, i.e., firms in which the government is the largest shareholder. Specifically, *%_Ind* has a positive and significant effect on

23

*ROA* and *ROE* in the subsample of firms in which the government is the top owner (columns (1) and (2)), but has no significant effect on either *ROA* or *ROE* in the subsample of firms in which non-state entities are top owners (columns (3) and (4)).

The principal role of the board of directors is to mitigate agency conflicts that arise in the corporate form of business due to the separation of ownership and control (Fama and Jensen, 1983). At the same time, firms with majority government ownership exhibit the ultimate separation between ownership (the nation's citizens) and control (managers/bureaucrats), and thus are subject to severe agency problems (see Megginson and Netter, 2001 for a survey of the relevant literature). In this context, our finding that board independence is positively related to firm performance in government-controlled firms suggests that monitoring by outside directors is more consequential in exactly those firms where agency costs are most severe.

It can also be argued, however, that the relevant governance issue in government-controlled firms is not an agency conflict due to a separation of ownership and control, but instead a conflict of interest between controlling shareholders (in this case the government) and minority shareholders as in La Porta, Lopez-de-Silanes, and Shleifer (1999). Prior literature suggests that controlling shareholders enjoy substantial private benefits of control often at the expense of minority shareholders (Nenova, 2003; Dyck and Zingales, 2004). One of the main differences between a private controlling shareholder versus the government as the controlling shareholder is that the latter may have incentives to pursue political and social objectives. Indeed, in a recent review article, Fan, Wei, and Xu (2011) note that government intervention through state ownership is a common feature of emerging markets. Chen, Sun, Tang, and Wu (2011) find that government ownership distorts firms' investment behavior and reduces investment efficiency. Therefore, our finding of a positive and significant relationship

24

between board independence and firm performance in government-controlled firms supports the view that increasing board independence can, at least in part, alleviate some of the inherent inefficiencies associated with government ownership. In the next part of our analysis, we explicitly investigate some of the potential channels through which independent directors can influence firm performance.

[Insert Table 6 here]

*5.3. Independent directors and tunneling*

The most often cited mechanism through which company insiders divert corporate resources is through related-party transactions (RPTs) (Dahya, Dimitrov, and McConnell, 2008). This type of diversion is often called "self-dealing" (Djankov, La Porta, Lopez-de-Silanes, and Shleifer, 2008) or "tunneling" (Johnson, La Porta, Lopez-de-Silanes, and Shleifer, 2000). Jiang, Lee, and Yue (2010) document that tunneling through intercorporate loans—a practice in which firm insiders divert corporate resources through loans to other entities (most of which unlisted) that are also under their control—is one of the most egregious examples of insider self-dealing in China.

Intercorporate loans appear on the balance sheets of Chinese listed firms as "other receivables," and prior to 2006 were almost never paid back. We follow Jiang, Lee, and Yue (2010) and estimate the magnitude of this brazen form of corporate abuse using ORECTA (other receivables scaled by total assets). To investigate the effect of board independence (*%_Ind*) on ORECTA, we estimate the following regression model, including a similar set of control variables as in Jiang, Lee, and Yue (2010) as well as firm and year fixed-effects:

$$ORECTA_{it} = \gamma \, \%\_Ind_{it} + \beta \, Control_{it} + d_i + d_t + \varepsilon_{it} \qquad \ldots \text{Eq. 4}$$

Jiang, Lee, and Yue (2010) report that tunneling through intercorporate loans was largely eradicated by 2006 after eight government ministries issued a joint statement stipulating that the top

25

management of the controlling entities would be held personally accountable through public disclosure and legal punishment if the intercorporate loans were not repaid. Therefore, we also estimate the relationship between *%_Ind* and ORECTA for two separate time periods, 1999 to 2005 and 2006 to 2012. We expect that any effect of board independence on ORECTA should be observed in the pre-2006 sample period. The results from estimating Equation (4) for our full sample period (1999-2012), sub-period 1999-2005, and sub-period 2006-2012 are presented in Table 7, Panels A, B, and C, respectively.

[Insert Table 7 here]

As reported in column (1) of Panel A in Table 7, *%_Ind* has a negative and significant coefficient, suggesting that board independence has a moderating effect on the tunneling of corporate resources through intercorporate loans in Chinese listed firms. A 10 percentage-point increase in the proportion of independent directors leads to a 0.17% decrease in intercorporate loans measured by ORECTA. For the average firm in our sample, this translates into about a 3.5 million RMB reduction in intercorporate loans.[12] When aggregating across all of the 2,057 unique firms in our sample, the total reduction in intercorporate loans amounts to about 7.3 billion RMB.

Columns (2) and (3) of Panel A in Table 7 examine the effect of *%_Ind* on ORECTA in firms where the top owners are state and non-state entities, respectively. The regression results indicate that the moderating effect of board independence on intercorporate loans is largely driven by government-controlled firms. Specifically, the estimated coefficient on *%_Ind* is -0.020 (*p*-value=0.03) in government-controlled firms and -0.017 (*p*-value=0.17) in non-government controlled firms. Thus, it appears that independent directors are more effective in constraining insider self-dealing in government-controlled firms where the incentives and opportunities for such activities are particularly high. This is also in line with our earlier findings, presented in Table 6, that the overall benefit

---

[12] The average firm in our sample has total book assets of 2,075 million RMB. A 10 percentage-point increase in the proportion of independent directors leads to: -0.0017*2,075 = 3.5 million RMB reduction in intercorporate loans.

26

associated with having more independent boards is significantly larger in government-controlled firms.

As expected, the results in Panels B and C of Table 7 suggest that the negative effect of board independence on intercorporate loans is only observed in the pre-2006 period, during which tunneling through intercorporate loans was rampant, but not in the period after 2006 when the central government took decisive action to curb this blatant form of corporate abuse.

### 5.4. Independent directors and investment efficiency

A firm's investment policy and its outcomes are central to the firm's long-term survival and growth. In a perfect world free of market imperfections, a firm's investment behavior is solely a function of its investment opportunity set (Modigliani and Miller, 1958). However, the real world is full of market frictions, such as information asymmetry and agency conflicts, which prevent firms from making optimal investment decisions.[13] In their study of investment policy in Chinese listed firms, Chen, Sun, Tang, and Wu (2011) find that government intervention represents another friction that distorts the investment behavior of state-controlled firms and leads to investment inefficiency. In this part of our analyses we investigate whether the positive effect of board independence on firm performance can be attributed, at least in part, to the ability of independent directors to positively impact firm investment behavior. In other words we test whether there is a positive relationship between board independence and investment efficiency in Chinese listed firms.

We follow Chen, Sun, Tang, and Wu (2011) and measure investment efficiency as the sensitivity of investment expenditures (*Invest_Expenditure*) to investment opportunities (measured by the log of Tobin's Q, *Ln(Q)*) using the following regression model:

$$Invest\_Expenditure_{it} = \gamma_1 \%\_Ind_{it} * Ln(Q)_{it} + \gamma_2 Ln(Q)_{it} + \beta Control_{it} + d_i + d_t + \varepsilon_{it} \quad \dots \text{Eq. 5}$$

Our main variable of interest is the interaction term between board independence and Tobin's

---

[13] For a survey of the relevant literature, please see Stein (2003).

Q (*%_Ind_{it}*Ln(Q))*. We include a similar set of control variables as in Chen, Sun, Tang, and Wu (2011) as well as firm and year fixed-effects. The regression results are presented in Table 8.


[Insert Table 8 here]


As shown in column (1) of Table 8, using our full sample, the interaction term *%_Ind * Ln(Q)* has a positive and significant (at the 1% level) coefficient, indicating that the sensitivity of investment expenditures to investment opportunities is stronger in firms with a higher proportion of independent directors. We also find that the interaction between Tobin's Q and government ownership (*Ln(Q)*TopOwner_State*) is negative and significant, confirming the main finding in Chen, Sun, Tang, and Wu (2011) that investment efficiency is lower in government-controlled firms. Finally, the results in columns (2) and (3) illustrate that the positive relationship between board independence and firms' investment efficiency holds for both government-controlled and non-government-controlled firms.

Given that agency conflicts and government intervention are some of the main factors explaining the suboptimal investment behavior of Chinese firms, we interpret the findings in Table 8 as supporting the view that board independence can moderate agency problems and reduce government intervention.


## 6. Additional results

As discussed earlier, the CSRC mandated in 2001 that by June 30, 2003 Chinese listed firms should have at least one-third of their boards comprised of independent directors. Taking advantage of this exogenous shock, we used the DID method as one of our empirical tests to provide the portfolio of evidence that independent directors enhance firm performance in Chinese listed firms (see Section 4.1.2).[14] While it is a useful identification tool in tackling the endogeneity problem, the 2001 mandate also raises several interesting questions: how much of the positive performance effect that we find for

---

[14] Many prior studies also rely on regulatory shocks to identify the effect of board structure on firm performance (see, e.g., Black and Khanna, 2007, Dahya and McConnell, 2007, Black and Kim, 2012, among others).

28

board independence is attributable to the 2001 mandate? Left to their own devices, will Chinese listed firms embrace independent boards and will board independence matter? Notably, answers to these questions have important policy implications. More specifically, so far we find strong evidence that board independence enhances firm performance. If there is also evidence that absent of the 2001 Guidelines, Chinese listed firms would not have adopted independent boards or adopted them too slowly, then an argument can be made that external pressure, either from regulators or other external stakeholders, on Chinese listed firms to appoint independent directors is warranted.[15]

The 2001 mandate also introduces possible alternative explanations for our main result—the positive relationship between board independence and firm performance. Chinese listed firms may be merely window dressing their boards, adopting the minimum number of independent directors to meet the regulatory requirement. In other words, Chinese independent directors are perfunctory, and the positive link we uncovered between board independence and firm performance is an artifact of other relationships. For example, the positive independence-performance relationship could be driven by the negative relationship between board size and firm performance. More specifically, assume two firms. Firm A has five directors and Firm B has six directors. To meet the 2001 mandate, both firms need to have a minimum of two independent directors, resulting in 40% board independence for Firm A that has a smaller board and 33% board independence for Firm B that has a bigger board. It is well established in the literature that smaller boards correlate negatively with firm performance (Yermack, 1996; Eisenberg, Sundgren, and Wells, 1998). Alternatively, as Figure 1 shows, the time-series variation in the fraction of independent directors is largely driven by the pre- and post-regulation variations. Therefore, the positive independence-performance relationship could also be driven by a sub-period effect. We include board size and firm and year fixed effects in all the performance regressions, which should mitigate the confounding effects due to variations in board size and time periods. Nonetheless more tests could be done to further alleviate these concerns. To address these

---

[15] The global regulatory push for more independent boards (e.g., the 1992 U.K. Cadbury Report and the 2002 U.S. Sarbanes-Oxley Act) is premised on the belief that more independent boards improve firm performance and that firms either do not voluntarily adopt independent boards or adopt them too slowly.

29

questions and alternative explanations, we perform four tests in this section.[16]

In our first test, we partition the sample into pre- (1999 to 2001) and post- (2003-2012) regulation periods, and re-run our baseline regression (Equation 1). We exclude 2002 because it is a transition year. If we find a significant effect for board independence in the pre-regulation period, then the evidence suggests that independent directors matter even prior to the 2001 Guidelines. Such a finding also alleviates the concern that the positive independence-performance relationship is driven by an increase in average firm performance in the post-regulation period. As Table 9 shows, the estimated coefficient of board independence is significantly and positively related to both ROA and ROE in the pre-regulation period, consistent with the idea that the positive independence-performance relationship is not entirely driven by the regulatory mandate and alleviating the concern of a sub-period effect. As expected, *%_Ind* is significantly and positively related to firm performance in the post-regulation period, consistent with the idea that the 2001 mandate has teeth and is beneficial to firm performance. While the coefficient of *%_Ind* is larger in magnitude in the pre-regulation period than in the post-regulation period, a test for equality of the coefficient estimates of *%_Ind* across the regressions reveals no statistical significance.

[Insert Table 9 here]

In our second test, we examine whether firms that voluntarily embrace higher levels of board independence than mandated by the 2001 Guidelines exhibit better performance. We re-run Equation 1 replacing *%_Ind* with *Extra_Ind* using the firm years from the post-regulation period (2003-2012). *Extra_Ind* is the *actual* number of independent directors minus the minimum number mandated by the 2001 Guidelines (the *mandated* number). To give examples of the *mandated* number of independent directors, if a firm has a five-member board, the number of mandated independent directors is two, while a seven-member board has a mandated independent director number of three. Therefore, this test sheds light on whether Chinese listed firms go beyond the regulatory mandate in utilizing

---

[16] We thank an anonymous referee for raising these interesting questions and alternative explanations.

30

independent directors and whether they benefit from doing so. This test also provides further evidence for whether the positive relationship between board independence and firm performance is driven by differences in board size. As Table 10 shows, *Extra_Ind* is significantly and positively related to both ROA and ROE. Therefore, results in Table 10 are consistent with those in Table 9 and the idea that some Chinese listed firms understand the value proposition of board independence and are able to reap greater economic benefits by appointing more independent directors than the minimum mandated by the 2001 Guidelines.

[Insert Table 10 here]

In our third test, we re-run Equation 1, replacing *%_Ind* with two new measures of board independence: (1) *Ind>50%*, an indicator variable set to one (and zero otherwise) if the percent of independent directors is greater than 50 percent, and (2) *Ind34%_50%*, an indicator variable set to one (and zero otherwise) if the percent of independent directors is no more than 50 percent but greater than one-third. Therefore, the base group in the regression consists of firms with no more than one third of board independence. A significantly positive coefficient of *Ind>50%* and/or *Ind34%_50%* would support the notion that Chinese independent directors are not perfunctory and firms have economic incentives to appoint such directors. As Table 11 shows, both *Ind>50%* and *Ind34%_50%* enter the ROA and ROE regressions with a significantly positive sign.

[Insert Tables 11 here]

In our fourth test, we provide additional evidence for the performance effect of independent directors from a different perspective—the meeting attendance record of independent directors. We conjecture that if independent directors are perfunctory, then it should not matter whether they miss

31

board meetings.[17] As Table 12 illustrates, the percent of missed board meetings by independent directors is significantly and negatively related to firm performance.[18] This finding contributes to our overall portfolio of evidence that independent directors play an important role in Chinese corporate governance.

[Insert Table 12 here]

To summarize, the additional tests in Section 6 provide corroborating evidence for the results from Section 4 that independent directors matter for China's listed firms. Given that the large increase in board independence in China's listed firms is driven by regulatory mandates, our results also suggest that regulation is useful in propelling firms to install more stringent internal governance controls. Therefore, our paper adds to the literature, which shows that stronger internal governance such as more independent board enhances firm performance, particularly when the legal protection of shareholder rights is weak (Dahya, Dimitrov, and McConnell, 2008; Black and Kim, 2012). Our results are also consistent with the global movement that has been gathering pace in the past two decades. Around the world, regulators and governance activists have been pressuring publicly traded firms to adopt more stringent governance measures (e.g., the 1992 U.K. Cadbury Report, Aggarwal, Erel, Ferreira, and Matos, 2011). In this regard, our study is timely in that it provides important empirical support for this movement.

## 7. Conclusion

In this paper, we use a panel data set covering more than 16,000 firm-years to investigate the effect of board independence on the performance of Chinese listed firms. Using several econometric techniques specifically designed to address the inherent endogeneity in the independence-performance

---

[17] Adam and Ferreira (2009) study the effect of women directors on firm performance and governance, and find that meeting attendance by directors plays an important role.

[18] We lose approximately 9 percent of our sample in these regressions, because board meetings data did not become available until 2004.

32

relationship, we find that the degree of board independence is positively and significantly related to firm performance, especially in government-controlled firms and in firms with lower information acquisition and monitoring costs. We further document that this positive relationship can be attributed, at least in part, to the ability of independent directors to prevent insider self-dealing and to improve investment efficiency in Chinese listed firms.

A unique characteristic of the Chinese corporate landscape is that most listed firms are former state-owned enterprises (SOEs) and the government is still the largest shareholder in many of these companies. In this context, our findings suggest that the appointment of independent directors, who can effectively monitor firm management and can uphold the goal of shareholder wealth maximization, can at least partially alleviate some of the inherent inefficiencies associated with having the government as a dominant shareholder. Other distinctive characteristics of the Chinese institutional environment are the lack of adequate investor protections provided by the legal system, and the nascent role of alternative governance mechanisms such as the market for corporate control and shareholder activism. The fact that board independence is an effective mechanism for alleviating agency conflicts in this weak investor protection environment supports the view that there is a substitution effect between internal and external corporate governance mechanisms.

Our overall finding is consistent with the spirit of the CSRC's *2001 Guidelines for Introducing Independent Directors to the Board of Directors of Listed Companies* that required listed firms' boards to have at least one-third of independent directors. While we have provided robust and comprehensive evidence of the relationship between board independence and firm performance in China's listed firms, there are still many important research questions that remain unanswered. For example, although we have examined some of the potential channels through which independent directors can impact firm performance, there are invariably other channels that can be explored. Future research can also examine the influence of individual director characteristics such as director ethnicity, busyness, proximity to corporate headquarters, etc. on directors' incentives and abilities to adequately fulfill their monitoring and advising roles.

33

## References

Adams, R., Hermalin, B., Weisbach, M., 2010. The role of boards of directors in corporate governance: A conceptual framework and survey. *Journal of Economic Literature* 48, 58-107.

Adams, R.B., Ferreira, D., 2007. A theory of friendly boards. *Journal of Finance* 62(1), 217-250.

Adams, R.B., Ferreira, D., 2009. Women in the boardroom and their impact on governance and performance. *Journal of Financial Economics* 94, 291-309.

Aggarwal, R., Erel, I., Ferreira, M., Matos, P., 2011. Does Governance Travel Around the World? Evidence from Institutional Investors, Journal of Financial Economics, 100, 154-181.

Aggarwal, R., Erel, I., Stulz, R., Williamson, R., 2009. Differences in governance practices between U.S. and foreign firms: Measurement, causes, and consequences. *Review of Financial Studies* 22(8), 3131-3169.

Allen, F., Qian, J., Qian, M., 2005. Law, finance, and economic growth in China. *Journal of Financial Economics* 77, 57-116.

Asch, S. E., 1951. Effects of group pressure upon the modification and distortion of judgments. Guetzkow. H. (ed.) *Groups, Leadership, and Men*, 177-190. Carnegie Press: Pittsburgh.

Bertrand, M., Mullainathan, S., 2003. Enjoying the quiet life? Corporate governance and managerial preferences. *Journal of Political Economy* 111, 1043-1075.

Boone, A.L., Field, L.C., Karpoff, J.M., Raheja, C.G., 2007. The determinants of corporate board size and composition: An empirical analysis. *Journal of Financial Economics* 85(1), 66-101.

Black, B., Khanna, V., 2007. Can corporate governance reforms increase firm market values? Event study evidence from India. *Journal of Empirical Legal Studies* 4(4), 749–796.

Black, B., Kim, W., 2012. The effect of board structure on firm value: A multiple identification strategies approach using Korean data. *Journal of Financial Economics* 104, 203-226.

Bruno, V., Claessens, S., 2010. Corporate governance and regulation: Can there be too much of a good thing? *Journal of Financial Intermediation* 19, 461-482.

Broome, L.L, Conley, J.M., Krawiec, K.D., 2011. Does critical mass matter? Views from the boardroom. *Seattle University Law Review* 34, 1049-1080.

Chen, S., Sun, Z., Tang, S., Wu, D., 2011. Government intervention and investment efficiency: evidence from China. *Journal of Corporate Finance* 17, 259-271.

Choi, J.J., Park, S.W., Yoo, S.S., 2007. The Value of outside directors: Evidence from corporate governance reform in Korea. *Journal of Financial and Quantitative Analysis* 42, 941-962.

Cicero, D., Wintoki, M.B., Yang, T., 2013. How do firms adjust their board structures? *Journal of Corporate Finance* 23, 108-127.

Cremers, K.J.M., Litov, L.P., Sepe, S.M., 2013, Staggered boards and firm value, revisited. Working paper.

Djankov, S., La Porta, R., Lopez-de-Silanes, F., Shleifer, A., 2008. The law and economics of

34

self-dealing. *Journal of Financial Economics* 88, 430-465.

Duchin, R., Matsusaka, J.G., Ozbas, O., 2010. When are outside directors effective? *Journal of Financial Economics* 96(2), 195-214.

Dahya, J., Dimitrov, O., McConnell, J.J., 2008. Dominant shareholders, corporate boards, and corporate value: A cross-country analysis. *Journal of Financial Economics* 87, 73 -100.

Dahya, J., McConnell, J.J., 2005. Outside directors and corporate board decisions. *Journal of Corporate Finance* 11, 37-60.

Dahya, J., McConnell, J.J., 2007. Board composition, corporate performance, and the Cadbury committee recommendation. *Journal of Financial and Quantitative Analysis* 42, 535-564.

Dyck, A., Zingales, L., 2004. Private benefits of control: An international comparison. *Journal of Finance* 59, 537-600.

Duchin, R., Matsusaka, J. G., Ozbas, O., 2010. When are outside directors effective? *Journal of Financial Economics* 96, 195-214.

Fan, G., Wang, X., 2006. In: *The Report on the Relative Process of Marketization of Regions in China*. The Economic Science Press, Beijing (in Chinese).

Fan, J.P.H., Wei, K.C.J., Xu, X., 2011. Corporate finance and governance in emerging markets: A selective review and an agenda for future research. *Journal of Corporate Finance* 17, 207-214.

Fan, J., Wong, T.J., Zhang, T., 2007. Politically-connected CEOs, corporate governance, and post-IPO performance of China's newly partially privatized firms. *Journal of Financial Economics* 84, 330-357.

Fama, E.F., Jensen, M.C., 1983. Separation of ownership and control. *Journal of Law and Economics* 26(2), 01-325.

Ferreira, M. A., Matos, P., 2008. The colors of investors' money: The role of institutional investors around the world. *Journal of Financial Economics* 88, 499-533.

Fedaseyeu, V., Linck, J.S., Wagner, H.F., 2014. The determinants of director compensation. Working paper.

Giannetti, M., Liao, G., Yu, X. 2014, The brain gain of corporate boards: A natural experiment from China. *Journal of Finance*, Forthcoming.

Guo, D., Jiang, K., 2013. Venture capital investment and the performance of entrepreneurial firms: Evidence from China. *Journal of Corporate Finance* 22, 375-395.

Hermalin, B.E., Weisbach, M.S., 2003. Boards of directors as an endogenously determined institution: A survey of the economic literature. *Economic Policy Review* 9, 7-26.

Jiang, L., Kim, J.-B., Pang, L., 2011. Control-ownership wedge and investment sensitivity to stock price. *Journal of Banking and Finance* 35, 2856-2867.

35

Jiang, G., Lee, C., Yue, H., 2010. Tunneling through intercorporate loans: The China experience. *Journal of Financial Economics* 98, 1-20.

Johnson, S., La Porta, R., Lopez-de-Silanes, F., Shleifer, A., 2000. Tunneling. *The American Economic Review* 90(2), 22-27.

Klapper, L. F., Love, I., 2004. Corporate governance, investor protection and performance in emerging markets. *Journal of Corporate Finance* 10, 703-728.

Knyazeva, A., Knyazeva, D., Masulis, R., 2013. The supply of corporate directors and board independence. *Review of Financial Studies* 26, 1561-1605.

La Porta, R., Lopez-de-Silanes, F., Shleifer, A., 1999. Corporate ownership around the world. *Journal of Finance* 54, 471-517.

Linck, J.S., Netter, J.M., Yang, T., 2008. The determinants of board structure. *Journal of Financial Economics* 87(2), 308-328.

Liu, Y., Wei, Z., Xie, F., 2014, Do women directors improve firm performance in China? *Journal of Corporate Finance*, forthcoming.

Maug, E., 1997. Boards of directors and capital structure: Alternative forms of corporate restructuring. *Journal of Corporate Finance* 3, 113-139.

Mei J., Scheinkman J., Xiong W., 2009. Speculative trading and stock prices: An analysis of Chinese AB share premia. *Annals of Economics and Finance* 10, 225-255.

McCahery, J. A., Sautner, Z., Starks, L. T., 2010. Behind the scenes: The corporate governance preferences of institutional investors. Working Paper.

Megginson, W. L., Netter, J. M., 2001. From state to market: A survey of empirical studies on privatization. *Journal of Economic Literature* 39, 321-389.

Modigliani, F., Miller, M., 1958. The cost of capital, corporate finance and the theory of investment. *The American Economic Review* 48(3), 261-297.

Nenova, T., 2003. The value of corporate voting rights and control: A cross-country analysis. *Journal of Financial Economics* 68, 325-351.

Petersen, M.A., 2009. Estimating standard errors in finance panel data sets: Comparing approaches. *Review Financial Studies* 22, 435-480.

Raheja, C.G., 2005. Determinants of board size and composition: A theory of corporate boards. *Journal of Financial and Quantitative Analysis* 40(2), 283-306.

Semadeni, M., Withers, M., Certo, S. T., 2013. The perils of endogeneity and instrumental variables in strategy research: Understanding through simulations. *Strategic Management Journal*.

Shanghai Stock Exchange. *A report on the corporate governance of Chinese publicly listed companies: Board independence and effectiveness*. Fudan University Press, 2004.

Shen, S., Jia, J., 2005. Will the independent director institution work in China? *Loyola Los Angeles International and Comparative Law Review*, 223-248.

36

Stein, J., 2003. Agency, information and corporate investment. In: Constantinides, G.M., Harris, M., Stulz, R. (Eds.), Handbook of the Economics of Finance. North-Holland, Amsterdam, The Netherlands.

Sun, Q., Tong, W., 2003. China share issue privatization: the extent of its success. *Journal of Financial Economics* 70, 183-222.

Wang, W., 2014. Independent directors and corporate performance in China: A meta-empirical study. Working paper.

Wang, C., Xu, H., 2011. Government intervention in investment by Chinese listed companies that have diversified into tourism. Tourism Management 32(6), 1371-1380.

Wintoki, M.B., Linck, J.S., Netter, J.M., 2012. Endogeneity and the dynamics of internal corporate governance. *Journal of Financial Economics* 105(3), 581-606.

Yang, T., Zhao, S., 2014. CEO duality and firm performance: Evidence from an exogenous shock to the competitive environment. *Journal of Banking and Finance*, forthcoming.

Yeh, Y.-H., Woidtke, T., 2005. Commitment or entrenchment?: Controlling shareholders and board composition. *Journal of Banking and Finance* 29, 1857-1885.

Yermack, D., 1996. Higher market valuation of companies with a small board of directors. *Journal of Financial Economics* 40(2), 185-211.

37

**Figure 1. Time trends of board independence in China, 1999-2012**

This figure depicts the time trend of board independence for our sample firms, 2,057 Chinese companies listed on the Shanghai and Shenzhen Stock Exchanges from 1999 to 2012.



a) Percent of independent directors on board



b) Percent of firms with independent directors

38

**Table 1: Summary statistics**

This table reports the summary statistics of key variables. The sample consists of 16,999 firm years or 2,057 unique firms from 1999 to 2012. Panel A reports summary statistics on firm performance measured by return on assets (*ROA*) and return on equity (*ROE*). *ROA* (*ROE*) is operating income before extraordinary items divided by total assets (total equity). Both *ROA* and *ROE* are truncated at the top and bottom 1%. Panel B reports summary statistics on board independence. *%_Ind* is the percent of independent directors on the board. *Ind_d1*, *Ind_d2*, *Ind_d3*, and *Ind_d4* are dummy variables that equal one when there is one, two, three, or four independent directors, respectively, and zero otherwise. *Ind_d5* equals one when there are five or more independent directors and zero otherwise. Panel C reports summary statistics of control variables. Ownership variables include: a dummy variable that equals one when the top owner of a listed firm is the State or a state-owned enterprise (SOE) (*Topowner_State*), a dummy variable that equals one when the top owner of a listed firm is an individual (*Topowner_Individual*), and the percent of B-shares and H-shares issued by a firm (*%_Foreign*). Proxies for the monitoring costs of a firm are: annualized standard deviation of weekly stock price returns (*Volatility*) and the geometric mean of sales growth rate over the past three years (*Sales_Growth*). Reported under other control variables are: the natural log of the number of directors on the board (*Ln_BoardSize*), a dummy variable that takes the value of one if the Chief Executive Officer chairs the board (*Duality*), investment expenditure (*Invest_Expenditure*) calculated as net cash payments for fixed assets, intangible assets, and other long-term assets divided by total assets at the beginning of the year (Chen, Sun, Tang and Wu, 2011), the natural log of the total assets (*Ln_Assets*), total liabilities divided by total assets (*Leverage*), and the natural log of the number of years that a firm is listed on an exchange (*Ln_FirmAge*).

| Variable | Obs | Mean | Std Error | 25th Ptcl | Median | 75th Ptcl |
|---|---|---|---|---|---|---|
| **Panel A: Performance measures** | | | | | | |
| ROA | 16999 | 0.036 | 0.063 | 0.011 | 0.037 | 0.067 |
| ROE | 16999 | 0.057 | 0.213 | 0.018 | 0.060 | 0.130 |
| **Panel B:   Independent directors** | | | | | | |
| %_Ind | 16999 | 0.304 | 0.132 | 0.308 | 0.333 | 0.364 |
| Ind_d1 | 16999 | 0.008 | 0.088 | 0 | 0 | 0 |
| Ind_d2 | 16999 | 0.095 | 0.293 | 0 | 0 | 0 |
| Ind_d3 | 16999 | 0.505 | 0.500 | 0 | 1 | 1 |
| Ind_d4 | 16999 | 0.182 | 0.386 | 0 | 0 | 0 |
| Ind_d5 | 16999 | 0.083 | 0.276 | 0 | 0 | 0 |
| **Panel C:   Control variables** | | | | | | |
| **Ownership variables** | | | | | | |
| Topowner_State | 16999 | 0.626 | 0.484 | 0 | 1 | 1 |
| Topowner_Individual | 16999 | 0.072 | 0.259 | 0 | 0 | 0 |
| %_Foreign | 16999 | 0.026 | 0.088 | 0.000 | 0.000 | 0.000 |

39

**Monitoring costs**

| | | | | | | |
|---|---|---|---|---|---|---|
| Volatility | 16999 | 2.457 | 2.814 | 0.898 | 1.556 | 2.866 |
| Sales_Growth | 16999 | 0.210 | 0.441 | 0.025 | 0.149 | 0.304 |

**Other**

| | | | | | | |
|---|---|---|---|---|---|---|
| Board_Size | 16999 | 9.573 | 2.210 | 9 | 9 | 11 |
| Duality | 16999 | 0.162 | 0.368 | 0 | 0 | 0 |
| Invest_Expenditure | 16999 | 0.071 | 0.095 | 0.014 | 0.044 | 0.097 |
| Ln_Assets | 16999 | 21.453 | 1.137 | 20.680 | 21.299 | 22.053 |
| Leverage | 16999 | 0.477 | 0.196 | 0.336 | 0.487 | 0.622 |
| Ln_FirmAge | 16999 | 1.947 | 0.664 | 1.386 | 2.079 | 2.485 |

40

**Table 2: Correlation matrix**

This table reports the correlations matrix of the variables used in our econometric analyses. Correlation coefficients significant at the 10% level or better are in **bold**. Variable descriptions are given in Table 1.

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 ROA | 1.00 | | | | | | | | | | | | | | | | | | |
| 2 ROE | **0.71** | 1.00 | | | | | | | | | | | | | | | | | |
| 3 %_Ind | **0.01** | 0.00 | 1.00 | | | | | | | | | | | | | | | | |
| 4 Ind_d1 | **-0.02** | **-0.02** | **-0.12** | 1.00 | | | | | | | | | | | | | | | |
| 5 Ind_d2 | **-0.07** | **-0.05** | **-0.10** | **-0.03** | 1.00 | | | | | | | | | | | | | | |
| 6 Ind_d3 | **0.03** | 0.00 | **0.32** | **-0.09** | **-0.33** | 1.00 | | | | | | | | | | | | | |
| 7 Ind_d4 | 0.01 | **0.02** | **0.26** | **-0.04** | **-0.15** | **-0.48** | 1.00 | | | | | | | | | | | | |
| 8 Ind_d5 | 0.00 | 0.01 | **0.25** | **-0.03** | **-0.10** | **-0.30** | **-0.14** | 1.00 | | | | | | | | | | | |
| 9 Topowner_State | **-0.04** | -0.02 | **-0.18** | 0.02 | 0.03 | **-0.17** | 0.04 | **0.04** | 1.00 | | | | | | | | | | |
| 10 Topowner_Individual | **0.09** | **0.07** | **0.15** | -0.02 | **-0.05** | **0.13** | -0.01 | -0.02 | **-0.36** | 1.00 | | | | | | | | | |
| 11 %_Foreign | **-0.03** | **-0.03** | -0.01 | -0.01 | -0.01 | **-0.04** | **0.03** | 0.02 | **0.09** | **-0.08** | 1.00 | | | | | | | | |
| 12 Volatility | **0.23** | **0.14** | **0.07** | -0.00 | **-0.08** | **0.06** | 0.01 | 0.01 | **-0.13** | **0.23** | **-0.05** | 1.00 | | | | | | | |
| 13 Sales_Growth | **0.23** | **0.17** | **0.04** | -0.01 | -0.00 | **0.02** | 0.00 | 0.00 | **-0.03** | 0.01 | **-0.03** | **0.11** | 1.00 | | | | | | |
| 14 Invest_Expenditure | **0.21** | **0.17** | **0.03** | -0.01 | **-0.04** | 0.01 | 0.01 | **0.04** | **-0.04** | **0.11** | **-0.04** | **0.15** | **0.15** | 1.00 | | | | | |
| 15 Ln_BoardSize | 0.01 | 0.00 | **-0.06** | **-0.03** | **-0.24** | **-0.30** | **0.29** | **0.46** | **0.15** | **-0.09** | **0.05** | -0.03 | -0.02 | 0.03 | 1.00 | | | | |
| 16 Duality | 0.01 | -0.00 | **0.03** | -0.00 | -0.02 | **0.06** | **-0.05** | -0.02 | **-0.19** | **0.19** | -0.03 | **0.09** | -0.01 | 0.03 | **-0.09** | 1.00 | | | |
| 17 Ln_Assets | **0.16** | **0.14** | **0.23** | **-0.04** | **-0.11** | -0.03 | **0.16** | **0.21** | **0.17** | **-0.10** | **0.20** | **0.06** | **0.09** | **0.12** | **0.20** | **-0.10** | 1.00 | | |
| 18 Leverage | **-0.37** | **-0.30** | **0.07** | -0.01 | -0.02 | -0.01 | **0.05** | **0.06** | **0.12** | **-0.22** | **0.05** | **-0.13** | **0.08** | **-0.09** | **0.06** | **-0.10** | **0.30** | 1.00 | |
| 19 Ln_FirmAge | **-0.16** | **-0.13** | **0.25** | **-0.04** | -0.02 | **0.10** | **0.05** | **0.08** | **0.14** | **-0.34** | **0.15** | **-0.22** | -0.02 | **-0.24** | 0.02 | **-0.13** | **0.22** | **0.34** | 1.00 |

41

**Table 3: Effect of board independence on firm performance**

This table reports the regression results using the ordinary-least-squares with firm and year fixed effects (*FE*), the difference-in-differences (*DID*), the two-stage least squares with instrumental variables (*IV-2SLS*), and the dynamic generalized method of moments (*GMM*) methods. All regressions are estimated with robust standard errors clustered at the firm level. Variable descriptions are given in Table 1. The *p*-value of the *t*-statistic of each coefficient is shown in italics. ***, **, and * indicate statistical significance at the 1%, 5%, and 10% level, respectively.

| | ROA | | | | | | | | ROE | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | (1) | | (2) | | (3) | | (4) | | (5) | | (6) | | (7) | | (8) | |
| | FE | | DID | | 2SLS | | Arellano Bond | | FE | | DID | | 2SLS | | Arellano Bond | |
| %_Ind | 0.029 | *** | | | 0.178 | * | 0.026 | ** | 0.138 | *** | | | 0.712 | * | 0.043 | |
| | *0.003* | | | | *0.052* | | *0.049* | | *0.000* | | | | *0.054* | | *0.428* | |
| Treated*Post_Legislation | | | 0.014 | ** | | | | | | | 0.047 | ** | | | | |
| | | | *0.014* | | | | | | | | *0.039* | | | | | |
| Ln_BoardSize | -0.004 | | -0.004 | | 0.001 | | 0.015 | | -0.025 | * | -0.031 | * | -0.009 | | -0.008 | |
| | *0.317* | | *0.322* | | *0.862* | | *0.571* | | *0.066* | | *0.052* | | *0.597* | | *0.916* | |
| Duality | -0.002 | | -0.002 | | -0.002 | | 0.010 | | -0.008 | | -0.009 | | -0.010 | | 0.066 | |
| | *0.445* | | *0.490* | | *0.276* | | *0.672* | | *0.336* | | *0.361* | | *0.211* | | *0.375* | |
| Topowner_State | -0.015 | *** | -0.015 | *** | -0.014 | *** | 0.001 | | -0.038 | *** | -0.038 | *** | -0.036 | *** | 0.015 | |
| | *0.000* | | *0.000* | | *0.000* | | *0.884* | | *0.005* | | *0.007* | | *0.007* | | *0.619* | |
| Topowner_Individual | 0.017 | | -0.009 | | 0.017 | | 0.006 | | 0.046 | | -0.030 | | 0.045 | | -0.049 | |
| | *0.152* | | *0.571* | | *0.163* | | *0.825* | | *0.388* | | *0.698* | | *0.411* | | *0.463* | |
| %_Foreign | 0.007 | | 0.014 | | 0.000 | | -0.005 | | -0.008 | | 0.015 | | -0.034 | | -0.034 | |
| | *0.829* | | *0.685* | | *0.990* | | *0.913* | | *0.957* | | *0.918* | | *0.820* | | *0.818* | |
| Volatility | 0.002 | *** | 0.003 | *** | 0.002 | *** | 0.007 | *** | 0.003 | *** | 0.004 | *** | 0.003 | *** | 0.018 | *** |
| | *0.000* | | *0.000* | | *0.000* | | *0.000* | | *0.000* | | *0.000* | | *0.000* | | *0.000* | |
| Sales_Growth | 0.029 | *** | 0.029 | *** | 0.029 | *** | -0.003 | | 0.083 | *** | 0.087 | *** | 0.083 | *** | 0.077 | |
| | *0.000* | | *0.000* | | *0.000* | | *0.857* | | *0.000* | | *0.000* | | *0.000* | | *0.182* | |
| Invest_Expenditure | 0.054 | *** | 0.057 | *** | 0.054 | *** | 0.047 | | 0.161 | *** | 0.182 | *** | 0.161 | *** | 0.162 | |
| | *0.000* | | *0.000* | | *0.000* | | *0.599* | | *0.000* | | *0.000* | | *0.000* | | *0.512* | |
| Ln_Assets | 0.013 | *** | 0.013 | *** | 0.012 | *** | 0.000 | | 0.052 | *** | 0.057 | *** | 0.049 | *** | 0.009 | |
| | *0.000* | | *0.000* | | *0.000* | | *0.925* | | *0.000* | | *0.000* | | *0.000* | | *0.460* | |
| Leverage | -0.171 | *** | -0.171 | *** | -0.172 | *** | -0.015 | | -0.476 | *** | -0.501 | *** | -0.478 | *** | -0.130 | |
| | *0.000* | | *0.000* | | *0.000* | | *0.704* | | *0.000* | | *0.000* | | *0.000* | | *0.313* | |

42

| | (1) | | (2) | | (3) | | (4) | | (5) | | (6) | | (7) | | (8) | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ln_FirmAge | -0.016 | *** | -0.024 | *** | -0.015 | *** | 0.002 | | -0.062 | *** | -0.086 | *** | -0.058 | *** | 0.012 | |
| | 0.000 | | 0.000 | | 0.000 | | 0.597 | | 0.000 | | 0.000 | | 0.000 | | 0.407 | |
| L2.Dep | | | | | | | 0.283 | *** | | | | | | | 0.245 | ** |
| | | | | | | | 0.001 | | | | | | | | 0.036 | |
| Obs | 16,999 | | 12,604 | | 16,982 | | 7,453 | | 16,999 | | 12,604 | | 16,982 | | 7,445 | |
| $R^2$ | 0.24 | | 0.25 | | 0.22 | | | | 0.16 | | 0.16 | | 0.13 | | | |
| First-Stage $F$ Test Statistics | | | | | 30.70 | | | | | | | | 30.70 | | | |
| Over-Identification Test $p$ Value | | | | | 0.38 | | 0.61 | | | | | | 0.80 | | 0.50 | |
| AR(1) of First-Differenced Residuals | | | | | | | 0.00 | | | | | | | | 0.03 | |
| AR(2) of First-Differenced Residuals | | | | | | | 0.75 | | | | | | | | 0.35 | |

**Table 4: The number of independent directors and firm performance**

This table reports the regression results of the *FE* method by replacing *%_Ind* with a set of dummy variables, *Ind_d1*, *Ind_d2*, *Ind_d3*, *Ind_d4*, and *Ind_d5*. All regressions are estimated with robust standard errors clustered at the firm level. Variable descriptions are given in Table 1. The *p*-value of the *t*-statistic of each coefficient is shown in italics. ***, **, and * indicate statistical significance at the 1%, 5%, and 10% level, respectively.

|  | (1) | | (2) | |
|---|---|---|---|---|
|  | **ROA** | | **ROE** | |
| Ind_d1 | 0.006 | | -0.004 | |
|  | *0.235* | | *0.863* | |
| Ind_d2 | 0.007 | * | 0.036 | ** |
|  | *0.068* | | *0.023* | |
| Ind_d3 | 0.016 | *** | 0.054 | *** |
|  | *0.001* | | *0.003* | |
| Ind_d4 | 0.016 | *** | 0.065 | *** |
|  | *0.001* | | *0.001* | |
| Ind_d5 | 0.016 | *** | 0.068 | *** |
|  | *0.003* | | *0.001* | |
| Ln_BoardSize | -0.008 | * | -0.049 | *** |
|  | *0.064* | | *0.004* | |
| Duality | -0.001 | | -0.007 | |
|  | *0.460* | | *0.371* | |
| Topowner_State | -0.015 | *** | -0.038 | *** |
|  | *0.000* | | *0.004* | |
| Topowner_Individual | 0.017 | | 0.046 | |
|  | *0.155* | | *0.387* | |
| %_Foreign | 0.007 | | -0.008 | |
|  | *0.836* | | *0.955* | |
| Volatility | 0.002 | *** | 0.002 | *** |
|  | *0.000* | | *0.000* | |
| Sales_Growth | 0.029 | *** | 0.083 | *** |
|  | *0.000* | | *0.000* | |
| Invest_Expenditure | 0.054 | *** | 0.161 | *** |
|  | *0.000* | | *0.000* | |
| Ln_Assets | 0.013 | *** | 0.052 | *** |
|  | *0.000* | | *0.000* | |
| Leverage | -0.171 | *** | -0.476 | *** |
|  | *0.000* | | *0.000* | |
| Ln_FirmAge | -0.016 | *** | -0.063 | *** |
|  | *0.000* | | *0.000* | |
| Obs | 16,999 | | 16,999 | |
| $R^2$ | 0.24 | | 0.16 | |

44

**Table 5: Independent directors, monitoring costs, and firm performance**

This table reports the regression results from using the *FE* method to test the moderating role of monitoring costs in the relationship between board independence and firm performance. Interaction terms, *%_Ind*Volatility* and *%_Ind*Sales_Growth*, are added to Equation (1) to test the moderating role. All regressions are estimated with robust standard errors clustered at the firm level. Variable descriptions are given in Table 1. The *p*-value of the *t*-statistic of each coefficient is shown in italics. ***, **, and * indicate statistical significance at the 1%, 5%, and 10% level, respectively.

| | (1) | | (2) | |
| --- | --- | --- | --- | --- |
| | **ROA** | | **ROE** | |
| %_Ind | 0.036 | *** | 0.159 | *** |
| | *0.000* | | *0.000* | |
| %_Ind*Volatility | -0.008 | *** | -0.022 | *** |
| | *0.000* | | *0.002* | |
| %_Ind*Sales_Growth | -0.034 | ** | -0.184 | *** |
| | *0.012* | | *0.002* | |
| Ln_BoardSize | -0.003 | | -0.024 | * |
| | *0.394* | | *0.086* | |
| Duality | -0.001 | | -0.007 | |
| | *0.471* | | *0.349* | |
| Topowner_State | -0.015 | *** | -0.038 | *** |
| | *0.000* | | *0.005* | |
| Topowner_Individual | 0.017 | | 0.044 | |
| | *0.157* | | *0.407* | |
| %_Foreign | 0.010 | | 0.003 | |
| | *0.761* | | *0.986* | |
| Volatility | 0.002 | *** | 0.002 | *** |
| | *0.000* | | *0.000* | |
| Sales_Growth | 0.028 | *** | 0.079 | *** |
| | *0.000* | | *0.000* | |
| Invest_Expenditure | 0.053 | *** | 0.161 | *** |
| | *0.000* | | *0.000* | |
| Ln_Assets | 0.014 | *** | 0.057 | *** |
| | *0.000* | | *0.000* | |
| Leverage | -0.173 | *** | -0.483 | *** |
| | *0.000* | | *0.000* | |
| Ln_FirmAge | -0.016 | *** | -0.062 | *** |
| | *0.000* | | *0.000* | |
| Obs | 16,999 | | 16,999 | |
| $R^2$ | 0.25 | | 0.16 | |

45

**Table 6: Independent directors, ownership structure, and firm performance**

This table reports regression results of the *FE* method for two subsamples. One subsample is composed of firms in which the government is the largest shareholder (*Topowner_State*=1). The other subsample is composed of firms with *Topowner_State* set to zero. All regressions are estimated with robust standard errors clustered at firm level. Variable descriptions are given in Table 1. The *p*-value of the *t*-statistic of each coefficient is shown in italics. ***, **, and * indicate statistical significance at the 1%, 5%, and 10% level, respectively.

| | (1) | | (2) | | (3) | | (4) | |
|---|---|---|---|---|---|---|---|---|
| | **ROA** | | **ROE** | | **ROA** | | **ROE** | |
| | Topowner_State=1 | | | | Topowner_State=0 | | | |
| %_Ind | 0.040 | *** | 0.142 | *** | 0.005 | | 0.102 | |
| | *0.000* | | *0.001* | | *0.796* | | *0.170* | |
| Ln_BoardSize | -0.002 | | -0.011 | | -0.004 | | -0.032 | |
| | *0.684* | | *0.473* | | *0.462* | | *0.243* | |
| Duality | -0.004 | | -0.015 | | -0.001 | | -0.007 | |
| | *0.176* | | *0.125* | | *0.866* | | *0.589* | |
| %_Foreign | 0.060 | | 0.136 | | 0.030 | | 0.257 | ** |
| | *0.158* | | *0.211* | | *0.527* | | *0.045* | |
| Volatility | 0.003 | *** | 0.003 | *** | 0.001 | *** | 0.002 | *** |
| | *0.000* | | *0.000* | | *0.000* | | *0.008* | |
| Sales_Growth | 0.033 | *** | 0.092 | *** | 0.024 | *** | 0.065 | *** |
| | *0.000* | | *0.000* | | *0.000* | | *0.000* | |
| Invest_Expenditure | 0.064 | *** | 0.185 | *** | 0.035 | *** | 0.116 | *** |
| | *0.000* | | *0.000* | | *0.000* | | *0.001* | |
| Ln_Assets | 0.014 | *** | 0.055 | *** | 0.012 | *** | 0.049 | *** |
| | *0.000* | | *0.000* | | *0.000* | | *0.000* | |
| Leverage | -0.192 | *** | -0.513 | *** | -0.141 | *** | -0.400 | *** |
| | *0.000* | | *0.000* | | *0.000* | | *0.000* | |
| Ln_FirmAge | -0.020 | *** | -0.065 | *** | -0.013 | *** | -0.062 | *** |
| | *0.000* | | *0.000* | | *0.001* | | *0.000* | |
| | | | | | | | | |
| Obs | 10,645 | | 10,645 | | 6,354 | | 6,354 | |
| $R^2$ | 0.265 | | 0.167 | | 0.202 | | 0.124 | |

46

**Table 7: Independent directors and tunneling**

This table reports regression results of Equation (4) using the FE method. The magnitude of corporate tunneling is estimated by ORECTA (other receivables scaled by total assets). *%_Ind*, two board variables (*Ln_BoardSize* and *Duality*), a similar set of control variables as in Jiang, Lee, and Yue (2010) as well as firm and year fixed-effects are used to explain ORECTA. *Block* is the percent of shares controlled by the top shareholder, and *Marketization* is an index measuring the development of the regional market (Fan and Wang, 2006). The dataset is further separated into two time periods, 1999 to 2005 and 2006 to 2012, as tunneling through intercorporate loans was largely eradicated by 2006 (Jiang, Lee, and Yue, 2010). The *p*-value of the *t*-statistic of each coefficient is shown in italics. ***, **, and * indicate statistical significance at the 1%, 5%, and 10% level, respectively.

Panel A: All years 1999-2012

| | Dependent Var: ORECTA Panel A: All Years | | | | | |
|---|---|---|---|---|---|---|
| | (1) All | | (2) TopOwner_State=1 | | (3) TopOwner_State=0 | |
| %_Ind | -0.017 | ** | -0.020 | ** | -0.017 | |
| | *0.020* | | *0.031* | | *0.168* | |
| Ln_BoardSize | -0.010 | *** | -0.011 | *** | -0.004 | |
| | *0.000* | | *0.000* | | *0.230* | |
| Duality | 0.000 | | 0.003 | | -0.001 | |
| | *0.790* | | *0.138* | | *0.711* | |
| L.ROA | -0.247 | *** | -0.239 | *** | -0.238 | *** |
| | *0.000* | | *0.000* | | *0.000* | |
| Block | -0.034 | *** | -0.036 | *** | -0.017 | *** |
| | *0.000* | | *0.000* | | *0.001* | |
| Ln_Assets | -0.004 | *** | -0.003 | *** | -0.006 | *** |
| | *0.000* | | *0.000* | | *0.000* | |
| TopOwner_State | -0.007 | *** | | | | |
| | *0.000* | | | | | |
| Marketization | -0.001 | *** | 0.000 | | -0.002 | *** |
| | *0.000* | | *0.259* | | *0.000* | |
| Obs | 16,815 | | 10,246 | | 6,569 | |
| $R^2$ | 0.26 | | 0.26 | | 0.32 | |

47

Panel B: Years 1999-2005

| | Dependent Var:ORECTA | | | | | |
|---|---|---|---|---|---|---|
| | Panel B: 1999-2005 | | | | | |
| | (1) | | (2) | | (3) | |
| | **All** | | **TopOwner_State=1** | | **TopOwner_State=0** | |
| %_Ind | -0.036 | ** | -0.034 | ** | -0.050 | |
| | *0.014* | | *0.036* | | *0.138* | |
| Ln_BoardSize | -0.012 | *** | -0.012 | ** | -0.010 | |
| | *0.005* | | *0.011* | | *0.318* | |
| Duality | 0.003 | | 0.003 | | 0.007 | |
| | *0.404* | | *0.455* | | *0.233* | |
| L.ROA | -0.437 | *** | -0.407 | *** | -0.474 | *** |
| | *0.000* | | *0.000* | | *0.000* | |
| Block | -0.051 | *** | -0.045 | *** | -0.070 | *** |
| | *0.000* | | *0.000* | | *0.000* | |
| Ln_Assets | -0.010 | *** | -0.008 | *** | -0.014 | *** |
| | *0.000* | | *0.000* | | *0.000* | |
| TopOwner_State | -0.008 | *** | | | | |
| | *0.001* | | | | | |
| Marketization | -0.002 | *** | -0.001 | | -0.006 | *** |
| | *0.001* | | *0.401* | | *0.000* | |
| | | | | | | |
| Obs | 6,158 | | 4,558 | | 1,600 | |
| $R^2$ | 0.24 | | 0.24 | | 0.31 | |

48

697

Panel C: Years 2006-2012

| | Dependent Var: ORECTA | | |
| | Panel C: 2006-2012 | | |
| | (1) | (2) | (3) |
| | **All** | **TopOwner_State=1** | **TopOwner_State=0** |
| %_Ind | 0.007 | 0.000 | 0.012 |
| | *0.327* | *0.970* | *0.260* |
| Ln_BoardSize | -0.007 *** | -0.006 ** | -0.005 |
| | *0.000* | *0.012* | *0.117* |
| Duality | -0.001 | 0.003 | -0.002 |
| | *0.415* | *0.132* | *0.174* |
| L.ROA | -0.138 *** | -0.110 *** | -0.153 *** |
| | *0.000* | *0.000* | *0.000* |
| Block | -0.020 *** | -0.024 *** | -0.012 *** |
| | *0.000* | *0.000* | *0.004* |
| Ln_Assets | -0.003 *** | -0.002 *** | -0.004 *** |
| | *0.000* | *0.000* | *0.000* |
| TopOwner_State | -0.004 *** | | |
| | *0.000* | | |
| Marketization | -0.001 *** | 0.000 | -0.002 *** |
| | *0.000* | *0.222* | *0.000* |
| | | | |
| Obs | 10,657 | 5,688 | 4,969 |
| $R^2$ | 0.18 | 0.19 | 0.23 |

49

**Table 8: Independent directors and investment efficiency**

This table reports regression results of Equation (5) using the FE method. Investment efficiency is measured as the sensitivity of investment expenditures to investment opportunities (measured by Tobin's Q). The main variable of interest is the interaction term between board independence and Tobin's Q (*%_Ind * Ln(Q)*). Also included in the right-hand side of the equation are *%_Ind*, two board variables (*Ln_BoardSize* and *Duality*), a similar set of control variables as in Chen, Sun, Tang, and Wu (2011), and firm and year fixed-effects. *CFO* is cash flow from operations. The *p*-value of the *t*-statistic of each coefficient is shown in italics. ***, **, and * indicate statistical significance at the 1%, 5%, and 10% level, respectively.

| | Dependent Var: Invest_Expenditure | | | | | |
|---|---|---|---|---|---|---|
| | (1) | | (2) | | (3) | |
| | All | | TopOwner_State=1 | | TopOwner_State=0 | |
| %_Ind | 0.036 | *** | 0.033 | *** | 0.015 | |
| | *0.000* | | *0.000* | | *0.244* | |
| %_Ind*Ln(Q) | 0.039 | *** | 0.032 | *** | 0.043 | *** |
| | *0.000* | | *0.002* | | *0.014* | |
| Ln(Q) | 0.013 | *** | 0.008 | *** | 0.006 | *** |
| | *0.000* | | *0.000* | | *0.005* | |
| TopOwner_State | -0.010 | *** | | | | |
| | *0.000* | | | | | |
| TopOwner_State*Ln(Q) | -0.007 | *** | | | | |
| | *0.003* | | | | | |
| Ln_BoardSize | -0.006 | * | *-0.006* | | *-0.010* | |
| | *0.067* | | *0.145* | | *0.109* | |
| Duality | *-0.002* | | *-0.001* | | *-0.001* | |
| | *0.299* | | *0.620* | | *0.757* | |
| CFO | 0.077 | *** | 0.080 | *** | 0.066 | *** |
| | *0.000* | | *0.000* | | *0.000* | |
| Leverage | -0.015 | *** | 0.001 | | -0.027 | *** |
| | *0.002* | | *0.920* | | *0.001* | |
| Ln_Assets | 0.021 | *** | 0.020 | *** | 0.026 | *** |
| | *0.000* | | *0.000* | | *0.000* | |
| Ln_FirmAge | -0.044 | *** | -0.038 | *** | -0.057 | *** |
| | *0.000* | | *0.000* | | *0.000* | |
| | | | | | | |
| Obs | 15,845 | | 9,730 | | 6,115 | |
| $R^2$ | 0.07 | | 0.05 | | 0.10 | |

50

699

**Table 9: Effect of independent directors on firm performance in pre- and post- regulation periods**

This table reports the effect of independent directors on firm performance, using the *FE* method, in pre- and post- regulation periods. Panel A reports regression results for the pre-regulation period (1999-2001). Panel B reports regression results for the post-regulation period (2003-2012). All regressions are estimated with robust standard errors clustered at firm level. Variable descriptions are given in Table 1. The *p*-value of the *t*-statistic of each coefficient is shown in italics. ***, **, and * indicate statistical significance at the 1%, 5%, and 10% level, respectively.

| | Panel A: Pre-Regulation Period 1999-2001 | | | | Panel B: Post-Regulation Period 2003-2012 | | | |
|---|---|---|---|---|---|---|---|---|
| | ROA | | ROE | | ROA | | ROE | |
| %_Ind | 0.033 | * | 0.264 | *** | 0.030 | ** | 0.121 | ** |
| | *0.079* | | *0.003* | | *0.016* | | *0.012* | |
| Ln_BoardSize | -0.009 | | -0.086 | * | -0.006 | | -0.033 | ** |
| | *0.342* | | *0.068* | | *0.228* | | *0.048* | |
| Duality | -0.005 | | -0.021 | | -0.003 | | -0.006 | |
| | *0.346* | | *0.405* | | *0.194* | | *0.513* | |
| TopOwner_State | -0.010 | | 0.001 | | -0.016 | *** | -0.043 | *** |
| | *0.502* | | *0.983* | | *0.001* | | *0.009* | |
| Topowner_Individual | | | | | 0.017 | | 0.046 | |
| | | | | | *0.135* | | *0.341* | |
| %_Foreign | -0.051 | | 0.260 | | 0.058 | | 0.042 | |
| | *0.543* | | *0.143* | | *0.173* | | *0.811* | |
| Volatility | 0.002 | ** | 0.003 | | 0.001 | *** | 0.002 | *** |
| | *0.025* | | *0.271* | | *0.000* | | *0.001* | |
| Sales_Growth | 0.027 | *** | 0.073 | *** | 0.026 | *** | 0.065 | *** |
| | *0.000* | | *0.003* | | *0.000* | | *0.000* | |
| Invest_Expenditure | 0.023 | * | 0.110 | * | 0.049 | *** | 0.127 | *** |
| | *0.089* | | *0.058* | | *0.000* | | *0.000* | |
| Ln_Assets | 0.017 | * | 0.051 | | 0.011 | *** | 0.047 | *** |
| | *0.069* | | *0.358* | | *0.000* | | *0.000* | |
| Leverage | -0.127 | *** | -0.475 | *** | -0.180 | *** | -0.446 | *** |
| | *0.000* | | *0.000* | | *0.000* | | *0.000* | |

51

| Ln_FirmAge | -0.053 | *** | -0.021 | | -0.010 | *** | -0.046 | *** |
|---|---|---|---|---|---|---|---|---|
| | *0.000* | | *0.739* | | *0.000* | | *0.000* | |
| Obs | 2,473 | | 2,473 | | 13,523 | | 13,523 | |
| $R^2$ | 0.19 | | 0.11 | | 0.22 | | 0.13 | |

52

**Table 10: The number of extra independent directors and firm performance**

This table reports the effect of extra number of independent directors on firm performance, using the *FE* method, in the post-regulation period (2003-2012). *Extra_Ind* is defined as the number of independent directors minus the number of mandated independent directors. All regressions are estimated with robust standard errors clustered at the firm level. Variable descriptions are given in Table 1. The *p*-value of the *t*-statistic of each coefficient is shown in italics. ***, **, and * indicate statistical significance at the 1%, 5%, and 10% level, respectively.

| | ROA | | ROE | |
|---|---|---|---|---|
| Extra_Ind | 0.003 | *** | 0.009 | ** |
| | *0.004* | | *0.025* | |
| Ln_BoardSize | -0.007 | | -0.038 | * |
| | *0.146* | | *0.028* | |
| Duality | -0.003 | | -0.006 | |
| | *0.206* | | *0.537* | |
| Topowner_State | -0.015 | *** | -0.043 | *** |
| | *0.001* | | *0.010* | |
| Topowner_Individual | 0.017 | | 0.046 | |
| | *0.134* | | *0.337* | |
| %_Foreign | 0.057 | | 0.041 | |
| | *0.179* | | *0.813* | |
| Volatility | 0.001 | *** | 0.002 | *** |
| | *0.000* | | *0.001* | |
| Sales_Growth | 0.026 | *** | 0.065 | *** |
| | *0.000* | | *0.000* | |
| Invest_Expenditure | 0.049 | *** | 0.126 | *** |
| | *0.000* | | *0.000* | |
| Ln_Assets | 0.011 | *** | 0.047 | *** |
| | *0.000* | | *0.000* | |
| Leverage | -0.180 | *** | -0.446 | *** |
| | *0.000* | | *0.000* | |
| Ln_FirmAge | -0.010 | *** | -0.047 | *** |
| | *0.000* | | *0.000* | |
| Obs | 13,523 | | 13,523 | |
| $R^2$ | 0.22 | | 0.13 | |

53

**Table 11: Majority independent boards and firm performance**

This table reports the effect of *Ind34%_50%* and *Ind>50%* on firm performance, using the *FE* method, in the post-regulation period (2003-2012). *Ind34%_50%* equals one if %_Ind>1/3 and %_Ind<=50% and zero otherwise. *Ind >50%* equals one if %_Ind>50% and zero otherwise. All regressions are estimated with robust standard errors clustered at the firm level. Variable descriptions are given in Table 1. The *p*-value of the *t*-statistic of each coefficient is shown in italics. \*\*\*, \*\*, and \* indicate statistical significance at the 1%, 5%, and 10% level, respectively.

| | ROA | | ROE | |
|---|---|---|---|---|
| Ind 34%_50% | 0.009 | *** | 0.032 | *** |
| | *0.002* | | *0.009* | |
| Ind >50% | 0.009 | * | 0.033 | ** |
| | *0.005* | | *0.013* | |
| Ln_BoardSize | -0.004 | | -0.028 | * |
| | *0.005* | | *0.017* | |
| Duality | -0.003 | | -0.005 | |
| | *0.002* | | *0.009* | |
| Topowner_State | -0.015 | *** | -0.043 | ** |
| | *0.005* | | *0.017* | |
| Topowner_Individual | 0.017 | | 0.046 | |
| | *0.012* | | *0.048* | |
| %_Foreign | 0.059 | | 0.046 | |
| | *0.043* | | *0.176* | |
| Volatility | 0.001 | *** | 0.002 | *** |
| | *0.000* | | *0.001* | |
| Sales_Growth | 0.026 | *** | 0.065 | *** |
| | *0.002* | | *0.007* | |
| Invest_Expenditure | 0.048 | *** | 0.124 | *** |
| | *0.007* | | *0.019* | |
| Ln_Assets | 0.011 | *** | 0.047 | *** |
| | *0.002* | | *0.007* | |
| Leverage | -0.180 | *** | -0.447 | *** |
| | *0.008* | | *0.031* | |
| Ln_FirmAge | -0.010 | *** | -0.046 | *** |
| | *0.003* | | *0.008* | |
| Obs | 13,523 | | 13,523 | |
| $R^2$ | 0.22 | | 0.13 | |

54

**Table 12: Missed board meetings by independent directors and firm performance**

This table reports the effect of *%_BoardMeetingMissed* on firm performance, using the *FE* method. *%_BoardMeetingMissed* is defined as the average number of board meetings missed by independent directors divided by the total number of board meetings. Data on *%_BoardMeetingMissed* is only available for the period of 2004 to 2012. All regressions are estimated with robust standard errors clustered at the firm level. Variable descriptions are given in Table 1. The *p*-value of the *t*-statistic of each coefficient is shown in italics. ***, **, and * indicate statistical significance at the 1%, 5%, and 10% level, respectively.

| | ROA | | ROE | |
|---|---|---|---|---|
| %_Ind | 0.030 | ** | 0.121 | ** |
| | *0.013* | | *0.051* | |
| %_BoardMeetingMissed | -0.026 | *** | -0.110 | *** |
| | *0.008* | | *0.038* | |
| Ln_BoardSize | -0.006 | | -0.035 | ** |
| | *0.005* | | *0.015* | |
| Duality | -0.002 | | 0.000 | |
| | *0.003* | | *0.009* | |
| Topowner_State | -0.012 | ** | -0.036 | * |
| | *0.005* | | *0.020* | |
| Topowner_Individual | 0.021 | * | 0.052 | |
| | *0.012* | | *0.051* | |
| %_Foreign | 0.055 | | 0.053 | |
| | *0.047* | | *0.182* | |
| Volatility | 0.002 | *** | 0.002 | *** |
| | *0.000* | | *0.001* | |
| Sales_Growth | 0.025 | *** | 0.059 | *** |
| | *0.002* | | *0.008* | |
| Invest_Expenditure | 0.046 | *** | 0.126 | *** |
| | *0.007* | | *0.018* | |
| Ln_Assets | 0.009 | *** | 0.043 | *** |
| | *0.002* | | *0.008* | |
| Leverage | -0.187 | *** | -0.447 | *** |
| | *0.008* | | *0.031* | |
| Ln_FirmAge | -0.012 | *** | -0.047 | *** |
| | *0.003* | | *0.008* | |
| Obs | 12,253 | | 12,253 | |
| $R^2$ | 0.22 | | 0.13 | |

55

**Appendix 1: Comparison of Independent Directors between China and the U.S.**

| | China | U.S. |
|---|---|---|
| | 1999-2012 | |
| Law requirement on %_Ind [a,b] | According to the Guidelines for Introducing Independent Directors to the Board of Directors of Listed Companies (the 2001 Guidelines), the board of directors of Shanghai- and Shenzhen- listed firms must have at least one-third independent directors. | According to the 2002 Sarbanes-Oxley Act (SOX), the board of directors of NYSE- and NASDAQ- listed firms must have a majority of independent directors. |
| Nomination process [a,b] | Independent directors can be nominated by a shareholder with more than 1% total shares. The nomination should be approved by the shareholder meetings. | Independent directors are nominated by the nominating committee or the board if the firm does not have a nominating committee. Post SOX, the nominating committee must consist entirely of independent directors. |
| Terms of appointment [a,b] | Elected every three years; maximum two terms; cannot be replaced during the term. | Elected annually; or in the case of staggered boards, a fraction of the directors are up for election each year. |
| Mean %_Ind [c,d] | 30% | 71% |
| Mean age of independent directors [c,d] | 52 | 62 |
| Percent of independent directors with firm shares [c,d] | 3.9% | 95.2% |
| Average annual compensation [c,d] | $7,607 | $178,320 |

56

Notes:

a. Chinese law requirements on %_Ind, nomination process, and the terms of appointment come from the Guidelines for Introducing Independent Directors to the Board of Directors of Listed Companies (the 2001 Guidelines), which was issued by CSRC in 2001. In reality, minority shareholders rarely have the opportunity to nominate independent directors. The Shanghai Stock Exchange (2004) reports that 70% of independent directors are nominated by a firm's top shareholders.

b. The U.S. law requirements on %_Ind and nomination process come from Linck, Netter, and Yang (2009). The U.S. practice in terms of director appointment comes from Cremers, Litov, and Sepe (2013). According to Cremers et al., about 60% of U.S. firms have staggered boards during 1999-2006. The ratio dropped to 45% in 2011.

c. Chinese information for %_Ind, director age, the percent of independent directors with firm shares, and annual compensation comes from the dataset we examined in this study.

d. The U.S. information for %_Ind, director age, and the percent of independent directors with firm shares comes from RiskMetrics, which covers S&P1500 firms. Annual compensation data come from Fedaseyeu, Linck, Wagner (2014), who study S&P1500 firms for the sample period of 2006-2010. Before 2006, U.S. publicly traded companies were not required to disclose director compensation for each individual board member.

57

**Highlights**

58

- We examine the relation between board independence and firm performance in China.

- Board independence is positively related to firm operating performance in China

- The effect of board independence is stronger in government controlled firms.

- Independent directors limit insider self-dealing and improve investment efficiency.

# P

## Principal-Agent Theory of Organizations

Robin Gauld
Department of Preventive and Social Medicine, Dunedin School of Medicine, University of Otago, Dunedin, New Zealand

## Synonyms

Agency theory; Institutional economics; Public choice theory; Theory of contracts; Theory of incentives

## Definition

Theory of interaction between an agent and the principal for whom they act, the point being to structure incentives so that the agent will act to benefit the principal.

## Introduction

The principal-agent theory of organizations ("agency theory" from here on) encapsulates the idea that public sector performance can be improved if incentive-based contracts between different actors are implemented. Principals will be more likely to achieve their desired outcomes, while agents will have clarity around work programs and goals. Agency theory has had considerable influence on the theory and practice of public administration and policy since its emergence in the 1970s. It was particularly instrumental in many high-income developed countries through the 1980s and 1990s, with often radical public sector reforms resulting. Its legacy has endured, with many public sector organizational and policy designs continuing to be underpinned by concepts derived from the theory. Based on institutional economics, agency theory has, therefore, provided a powerful and all-encompassing framework for public sector organization. As such, there has been much written about agency theory itself, and about public sector contracting, which is a central tenet (see for instance Ashton et al. 2004; Klingner et al. 2002; Lane 2001; Pallesen 2004; Pinch and Patterson 2000). There is less literature that discusses the longer-term policy outcomes when agency theory has been an overarching influence on public sector organization (Gauld 2007).

This chapter provides an overview of agency theory, including the key ideas behind it and the organizational and policy arrangements that are derived from it. The chapter outlines the benefits that might be expected to result when the theory is applied to public sector organization. It also highlights agency theory's shortcomings. Finally, it notes areas where theoretical extension is demanded.

## What Is Agency Theory?

Agency theory has its foundations in two ideas which were developed through the study of the economics of organizational and institutional behavior.

First was the notion proposed by public choice theorists that self-interest is the primary motivation behind the activities and behavior of individuals and the organizations that they work for (Mitchell and Simmons 1994; Mueller 1989). Through this lens, people in public employment and public organizations are viewed as "rational utility maximizers," meaning that each seeks to advance their own interests as would a private business or private sector employee in pursuit of profit or an increased salary. In this way, government officials and organizations are presumed to be in pursuit of only budgetary expansion; politicians, for their part, are motivated by the prospect of an expanded share of votes; interest groups, who represent specific sector groups or services users, are only concerned with furthering their own ends and those of their members. Public choice theory suggests that the result of such behaviors is a state which is larger than it should be, along with policies which are designed primarily to serve voter preferences and boost voter support for politicians and political parties, and an economy which is skewed through meeting the demands of selected interests over and above those of the broader public interest. To counter this, public choice theorists advocate limits on the power of politicians, interest groups, and public officials, as well as the implementation of financial incentives and sanctions to ensure appropriate performances.

The second idea, again based on a presumption that self-interest and rational utility maximization drive behavior, is the view that all of life, including public work, private lives, and organizational activities, can be viewed as a set of relationships between different parties (Moe 1984; Perrow 1986). The details of these relationships, once defined, and the requirements of the various parties in any particular venture or activity can be itemized and then written into a formal contract. The contract can, in turn, be deployed for purposes such as setting expectations and objectives of contract partners and for establishing performance assessment and accountability expectations.

Advocates of agency theory and related organizational arrangements presume that contracting will align the interests of principals (those wanting something done, such as politicians, funding agencies, or chief executives) and their agents (e.g., government officials and organizations or non-government and private service providers of public services). The result, in theory, is that the achievement of principals' objectives will be maximized, resulting in a more efficient and effective policy and service delivery outcomes. Alongside of this, the self-interested behavior of agents will be stemmed and focused on principals' goals via various incentives and sanctions. These might include anything from withholding of service delivery payments through to organizational or individual employee performance bonuses (detailed discussion of agency and related theory can be found in Dixit 2002; Laffont and Mortimort 2002; Mueller 1989; Self 1993; Stretton and Orchard 1994; Wallis and Dollery 1999; Walsh 1995).

## Factors that Complicate Application of Agency Theory

Complicating agency theory is a series of behavioral factors encapsulated by the terms "adverse selection" and "moral hazard," as well as the very nature of the public sector and government. Adverse selection results from the existence in any relationship of what are called "information asymmetries." This refers to the simple fact that one party (for instance, an experienced and skilled public servant) may be likely to have more knowledge, and therefore be at an advantage, than another party (such as a politician). Adverse selection can occur when, for example, a principal (a politician) is not able to gain sufficient knowledge about an agent's (a public servant's) background, motivations, or capabilities prior to entering into a contractual relationship (Perrow 1986). It can pose particular problems in any

contracting situation and require considerable investigation to ensure that the potential for adverse selection related-difficulties is reduced.

The risk of "moral hazard" arises once a contract has been agreed to. It stems from the fact that, on a day-to-day basis, principals are not able to observe most agent activity (Moe 1984). Principals are, therefore, reliant on agents carrying out tasks and performing at a level as specified in a contract. The ever-present prospect of moral hazard means an ongoing requirement for monitoring. This can result in considerable costs to both contracting parties, as well as goal-displacement behavior on behalf of agents where they place a disproportionate emphasis on work that is specifically subject to monitoring. By goal displacement, this means that agents focus on monitored goals, to the detriment of other organizational and individual goals that may not be directly monitored.

A further factor complicating agency theory is a range of circumstances particular to government and public sector work. These include that policy refinement is frequently left to the implementation process and is routinely the responsibility of agents, being public officials and not principals (Hill and Hupe 2002); most government agencies have several and often conflicting tasks and objectives which can be difficult to define and itemize (Wilson 1989); multiple principals and agents characterize the public sector as do situations in which principals often double as agents; the public sector tends to lack competition, at least in terms of core government non-trading functions (Allison 1979); and the public workforce and agencies are motivated by a complex array of factors, only one of which might be financial incentives (Dixit 2002; Le Grand 2003; Thaler 2015; Wilson 1989).

## Improving Agency Theory

There is now considerable international experience with application of agency theory to public sector organization. It has provided the foundations for public sector reforms in a range of high-income countries. While there is an absence of research into whether long-term public sector performance has improved as a result, and it could be difficult to determine direct causality, it may be fair to suggest that performances have continued to be questionable. Indeed, politicians in most high-income countries, along with international agencies, continue to demand improvement in a context where policy challenges are increasingly complex (e.g., issues ranging from how to deal with population aging and chronic disease through to inequality and climate change which require considerable cross-government coordination and long-range planning). This is despite agency theory having an ongoing and influential role as a framework for establishing the roles of principals and agents and, in turn, for managing policy development and implementation.

The above discussion brings to the fore questions over the applicability of agency theory to complex and changing policy and management issues and points to several agency theory deficiencies. Most of these have been alluded to elsewhere but not incorporated explicitly into the theory (see, e.g., Dixit 2002; Laffont and Mortimort 2002; Le Grand 2003). In other words, they do not feature prominently in discussions by agency theorists about the implications of, and prospects for, arrangements derived from the theory. There are four key shortcomings. These include:

1. Principals (in this case, political leaders and public officials at different levels of government and the health system) often may not have sufficiently detailed knowledge of what they want when setting parameters and building incentives and goals for agents. If they do, and decide upon a certain policy path, then this has implications if political and policy preferences are subsequently altered. To counter such problems, principals may need a longer-term view, as discussed below.
2. Principals may not recognize in advance the ramifications of the directions they set. In response, it might be suggested that principals clearly need to model in detail the possible outcomes of various policy options. Such an approach, however, would be subject to the

widely noted limits of "rational" policymaking (Lindblom 1959).

3. Agency theory fails to adequately account for situations when principals regularly change. Such changes, for instance in political leadership, can be a harbinger for change in policy directions and the organization of the public sector. Agency theory may simply need to be amended to include reference to the fact that arrangements inspired by it cannot be relied upon to produce continuity in transitions between principals.

4. Following the previous point, when the administrative systems within which agents work are regularly restructured by principals, with the potential to induce confusion and chaos, this undermines the assumption inherent to agency theory that principals are in control and capable of providing consistent direction over time; it also disrupts administrative processes and continuity, creating obstacles to effective individual and aggregate agent activity. A consequence is that agents may develop their own methods and systems for working and fail to coordinate with one another. If incentives are to be relied upon for complex organization and in scenarios of transition, then their design needs to be sophisticated. A simple vertical contract between principals and agents may not be enough; horizontal contracts between agents may also be required, and these may need to be coordinated by principals. By implication, this means that principals need the detailed knowledge and foresight discussed in points 1 and 2 above, along with commitment to stability implied in point 3.

Thus, for agency theory to be an effective organizing principle now and into the future, more intricacy in its development may be required (Deacon 2004). Yet even if an approach that accounted for the four points above were applied, this may be too constrained for the exigencies of politics, the multitudinous motivations of individuals and public organizations, the ever-changing and increasingly complex nature of administration and society, and the challenges for the foreseeable future. It may also prove to be administratively demanding, raising questions over transaction costs.

If so, then alternatives may need to be considered. One possibility is for increasing the level of centralized policy control and sector oversight, in short, the growth of a more comprehensive iteration of agency theory that incorporates the coordination of principals and agents. If contracting is to remain a fundamental principle of public sector organization, then longer-term contracts, combined with standardized objectives and funding levels (inherent to comprehensive policy), may promote more commitment among agents to centrally driven directions. This points to the need for development of a consensus-based variant, bringing together principals and agents, with an aim of stewardship over a subset of contractual agreements between interested parties. In practice, this may mean lengthy consultation processes that involve policymakers, key interests, industry, and provider groups with a primary aim of forging long-term goals. Such an approach might be underpinned by a focus on outcome-oriented policies being increasingly pursued by various national governments (Baehler 2003; Christensen and Laegreid 2013; Hoque 2008). Forging a long-term, cross-government view may carry with it a risk of committing to particular directions that, following any future political change, could be rejected. Consensus may, also, over time lead to policy embeddedness and establishment of a set of institutions that would be resistant to future reform (see Blank and Burau 2004; Putnam 1993: 179; Wilsford 1994). This may be appropriate when effective responses are demanded to complex policy problems such as those listed earlier in this chapter.

## Conclusion

Short of the developments described in the previous section, which might be seen by agency theory proponents as undermining competitive incentive systems, the rudimentary method of simple contracts between principals and agents may fail to provide an effective developmental foundation for multifarious administrative issues.

## Cross-References

## References

Allison G (1979) Public and private management: are they fundamentally alike in all unimportant respects? In: Shafritz J, Hyde A (eds) Classics of public administration. Brooks/Cole Publishing Company, Pacific Grove

Ashton T, Cumming J, McLean J (2004) Contracting for health services in a public health system: the New Zealand experience. Health Policy 69(1):21–31

Baehler K (2003) "Managing for outcomes": accountability and thrust. Aust J Public Adm 62(4):23–34

Blank RH, Burau V (2004) Comparative health policy. Palgrave Macmillan, Houndmills

Christensen T, Laegreid P (eds) (2013) Transcending new public management: the transformation of public sector reforms. Ashgate, Surrey

Deacon A (2004) Book review: Julian Le Grand, "motivation, agency and public policy: of knights, knaves, pawns and queens". J Soc Policy 33(3):503–506

Dixit A (2002) Incentives and organizations in the public sector: an interpretive review. J Hum Resour 32(4):696–727

Gauld R (2007) Principal-agent theory and organisational change: lessons from New Zealand health information management. Policy Stud 28(1):17–34

Hill M, Hupe P (2002) Implementing public policy. Sage, London

Hoque Z (2008) Measuring and reporting public sector outputs/outcomes: exploratory evidence from Australia. Int J Public Sect Manag 21(5):468–493. doi:10.1108/09513550810885787

Klingner DE, Nalbandian J, Romzek BS (2002) Politics, administration, and markets – conflicting expectations and accountability. Am Rev Public Adm 32(2):117–144

Laffont J-J, Mortimort D (2002) The theory of incentives: the principal-agent model. Princeton University Press, Princeton

Lane JE (2001) From long-term to short-term contracting. Public Adm 79(1):29–47

Le Grand J (2003) Motivation, agency and public policy: of knights, knaves, pawns and queens. Oxford University Press, Oxford

Lindblom C (1959) The science of muddling through. Public Adm Rev 19(2):79–88

Mitchell W, Simmons R (1994) Beyond politics: markets, welfare, and the failure of bureaucracy. Westview Press, Boulder

Moe T (1984) The new economics of organization. Am J Polit Sci 28:739–775

Mueller D (1989) Public choice II. Cambridge University Press, Cambridge

Pallesen T (2004) A political perspective on contracting out: the politics of good times. Experiences from Danish local governments. Governance 17(4):573–587

Perrow C (1986) Complex organizations: a critical essay. Random House, New York

Pinch PL, Patterson A (2000) Public sector restructuring and regional development: the impact of compulsory competitive tendering in the UK. Reg Stud 34(3):265–275

Putnam R (1993) Making democracy work: civic traditions in modern Italy. Princeton University Press, Princeton

Self P (1993) Government by the market? The politics of public choice. Westview Press, Boulder

Stretton H, Orchard L (1994) Public goods, public enterprise, public choice: theoretical foundations of the contemporary attack on government. Macmillan Press, Houndmills

Thaler RH (2015) Misbehaving: the making of behavioural economics. Allen Lane, London

Wallis J, Dollery B (1999) Market failure, government failure, leadership and public policy. Macmillan, London

Walsh K (1995) Public services and market mechanisms: competition, contracting and the new public management. Macmillan, London

Wilsford D (1994) Path dependency, or why history makes it difficult but not impossible to reform health systems in a big way. J Public Policy 14(3):251–283

Wilson JQ (1989) Bureaucracy: what government agencies do and why they do it. Basic Books, New York

# After NPM, curb your enthusiasm for the Principal-Agent theory[1]

Sten Widmalm

**Abstract**

Today the failures of New-Public-Management-inspired ideas to address classic challenges to a public administration, and the way NPM subsequently created new dysfunctions in state apparatuses, inspire us to scrutinize and decide which theoretical components in this field of research deserve to be retained, and which should be abandoned. This is not motivated by any conviction that, somewhere out there, there is a new "grand theory" that simply needs to be discovered. It is argued that not only the future challenges in public administrations, but also those still current, require us to make use of the wide range of analytical tools already available. It also requires a stance against reductionist economic theory. To make this point, this article focuses on the Principal Agent Theory and its origins, underlining that it remains an increasingly popular approach to the analysis of public administration, but arguing that both normatively and theoretically this theory is more problematical than is usually recognized.

To understand what makes people do "the right thing" in public administration – even when no-one is looking – can be said to be a number one challenge from both a practical and a research perspective. However, the New Public Management (NPM) administration models that are still employed today in many parts of the world – in particular in Sweden – seem to take us further from working solutions. This article argues that it will be easier to find alternatives to NPM public administration strategies, and that public administration research endeavours will move forward more easily, by relying less on the Principal-Agent theory.

Sten Widmalm works in the Department of Government at Uppsala University.
E-mail: sten.widmalm@statsvet.uu.se

Criticism of NPM-inspired models in the public administration discourse has long been widespread in the theoretical literature. It is often said that the discourse "moved on" a long time ago, although in practice NPM has survived (Pollitt 2014). More recently, however, continuing criticism in the research literature has manifested itself more frequently in the real world. NPM is criticized in the public debate for providing perverse incentive structures, overwhelming and ineffective evaluation structures, a dehumanizing work environment, and for deprofessionalizing a number of occupations in the health, education, research and administration sectors (Ahlbäck Öberg 2010; Zaremba 2013; Liedman 2012; Brante et al., 2015; Ahlbäck Öberg and Bringselius, 2015).

However, the claim that the research discourses have "moved on" and that public administrations are mainly waiting to "catch-up" is far from unproblematic for many reasons.

To begin with, public administration has never been known for following the lead of researchers. At least not in any simple way. Politics and economics shape the conditions for reform (see the introductory chapter) and what is finally adopted as a new strategy may hinge more on influences from market forces than on articles published by researchers. Policy makers who are looking for new ways forward are, however, somewhat frustrated by the fact that leading researchers are not recommending a simple "new" model that will provide solutions to the public administration challenges now that New Public Management is no longer in vogue. But here the message should be that one of the most important lessons to be learned from studying the deficiencies in NPM and the application of NPM solutions is that it assumes a "one-size fits all" kind of thinking. In an increasingly complex society, more competent and purpose-designed solutions, combined with political ideas, seem to be the way forward. Designing a public administration is not only a realm for technocrats. In a democracy it is also the subject of politics. Naturally this is not to say that we do not need help from broader research perspectives in discussions on how public administration solutions can and should be designed. And for this reason self-examination in the research discourse is not only welcome but also necessary.

This self-examination, it will be argued here, contains several guiding insights that can provide wider perspectives which may then help researchers and public administration policy makers to move forward. Several important insights have already been presented and demonstrated in the other articles included in this issue of *Statsvetenskaplig tidskrift*. However, there is one implication from several of these contributions that needs to be spelled out. And that is that in some vital respects the field of public administration studies has not really "moved on" to the extent that has been claimed. The more policy makers and public administration researchers try to present a position outside the NPM paradigm, the more evident it becomes that at least some central aspects of the paradigm have taken on the role of a meta- or supra-

ideology. Tingsten wrote that democracy was once a "supra-ideology" (Herbert Tingsten, [1945] 1960: 57). This being so, "one is a democrat, but also conservative, liberal or a socialist." Tingsten then expressed his worries about democracy losing its position as a supra-ideology. Here the problem is the opposite. The need is rather for the NPM mindset to loosen its grip. In the realm of public administration today, the NPM mindset has such a hegemonic role. To paraphrase Tingsten, "one is an NPM public servant, but also a teacher, a health care worker or a bureaucrat".

One of the most important illustrations of the extent to which the public administration research discourse is bogged down like this is manifested in the prominent role held by the Principal Agent (PA) theory. This is firmly rooted in a Rational Choice perspective and it is used today in general public administration studies and more specifically in areas of governance and corruption studies. It generates advice on incentive structures for Public Administration policy advisors and strategists, but several of its basic assumptions can be considered burdened with many of the things that NPM was criticized for. It will be argued here that PA theory is not only weak as a theory but may even have some detrimental effects if applied in practice. The discussion here considers the descriptive, explanatory and normative qualities of PA theory. Some references will also be made to Collective Action theory, Social Capital studies, and studies in public administration culture to illustrate these points. The message in this article is not that the criticism of PA theory forces us to throw out rational choice theory (RC) with the bathwater. The message is rather that we should not let RC-inspired theories such as PA theory throw out everything else.

## From Rational Choice, to the prominence of the Principal Agent theory

In the US in the late 1980s public trust in the government was at an all-time low and only "one in four Americans expressed confidence in government to "do what is right"" (Perry and Recascino Wise 1990). The situation was quite similar in many places in the West and researchers focusing on public administration performance struggled with perspectives and theories that could capture what guided the actions of public servants. At the time the one perspective which was strongly on the rise politically and among researchers was the "public choice movement." The dominant idea it proposed was that individuals were mainly motivated by self-interest. It would be unfair to blame the original work of Anthony Downs, Mancur Olson and Kenneth Arrow for all the extreme theoretical interpretations which have followed from the Rational Choice models that were established in the 1950s and 1960s. However, Public Choice (PC), which passionately advocated that market principles be applied in public administration (and many other areas as well), followed immediately

(Buchanan and Gordon Tullock 1962), and in the realm of public administration this later evolved into New Public Management. The public choice advocates essentially asserted that economic incentives guide human action and therefore that self-interest would certainly interfere with public interest in public administration enterprises. A quite credible explanation was provided for why bureaucracies tend to grow and why self-interest tends to replace public interest. The solution was not to fight this force, but to endorse it and consequently to steer and design the organisational public administrations according to this understanding of the human psyche.

The contrast could not have been starker than with predecessors like Herbert Kaufman who emphasised socialization and mores in various forms to explain why public servants would do the right thing in situations that were almost always unmonitored (Kaufman 2006). His insightful study of forest rangers in America explained how norms and altruistic goals were successfully conveyed and then sustained within a relatively small part of the government which had the responsibility of protecting vast regions of the country with infinite values. Without micromanagement, detailed checklists, or obvious carrot-and-stick incentives, the forest ranger in America could almost always be expected to do the right thing – alone in the woods, with possibly just a bird or a bear watching.

What could have happened when Public Choice picked up the challenge to predecessors like Kaufman was that political forces and key policy makers could have decided to make use of the wide range of theories existing between the two end-points described here. However, the political force behind the ideological content associated with Public Choice theories provided for a one-sided application that was unprecedented. A paradigm shift followed (Kuhn 1970).

Fast forward twenty years and the two most obvious examples are the doctrines of government applied by Ronald Reagan in the US, and Margaret Thatcher in the United Kingdom. Fast forward another decade and the application of NPM spread in most advanced economies of the West – although with differences in pace and scope. Sweden, the United Kingdom and New Zealand have in common the fact that some of the most radical reforms geared toward the market-inspired models were introduced by Labour governments in the 1990s. Simultaneously, in the public administration discourses more warnings of the detrimental side-effects of applying simplistic models of governance to modern public administration challenges were issued (Power 1997; Dunleavy and Hood 1994; Lindgren 2014; Widmalm et al., 2014; Ahlbäck Öberg and Widmalm 2012; Ahlbäck Öberg and Widmalm 2013; Liedman 2012). However, it would take another decade before leading political forces would start to listen and even consider the warnings that were raised. This brings us to today when public administration researchers not only speak of the detrimental effects of NPM but also try to point out ways forward.

It is to a large extent true that the downside of applying NPM or a Public Choice paradigm in a one-sided way has been "dealt with" in the public administration literature. However, the argument here is that an exaggerated economist mode of thinking still provides something akin to a supra-ideology in the discourse on governance. And in its most problematic forms it hampers efforts to move forward analytically and to make use in practice of the variety of tools that we not only need, but also to a large extent already have, when tackling modern governance problems.

The focus here is on the prominent theoretical trajectory called Principal-Agent theory, which has been chosen since it is so often promoted as one of the most useful tools for investigating and understanding problems relating to corruption – when people definitely do the wrong thing when they are not being watched. The way of thinking advocated by the PA theorists crowds out (to borrow the vocabulary of the NPM critics) other fields of knowledge and it reproduces reductionist perspectives on human behaviour and what guides action – just like NPM. In particular it does so when it comes to what makes people do the right thing in public administration.

In the opening sentence of her seminal contribution "Corruption and Government" Susan Rose-Ackerman confidently proclaims that "economics is a powerful tool for the analysis of corruption" (Rose-Ackerman 1999). In order to sort out distinctions between gifts and bribes in relation to the concept of corruption, she makes use of PA theory. There is no doubt that Rose-Ackerman's analysis is an important contribution to corruption studies focusing on the former Eastern bloc. It would be easy to assume that it was the empirical knowledge which she provided, rather than the application of PA theory, which made her contribution a success. Nevertheless it seems that Rose-Ackerman's position on PA theory must have struck a chord for some.

PA perspectives became more common as NPM models gained ground. They evolved simultaneously, since NPM and the terminology dates back at least to the mid-1970s. However, in corruption studies we can see that it has been emphasized since the late 1980s (Jensen and Meckling 1976; Klitgaard 1988). The general idea comes from economists who regard actors as either customers or clients. In PA-theory today we see that the principal is an actor that may represent a certain interest – for example the public interest. The agent is another actor that is supposed to carry out some action or actions for the principal – for example provide some public services. However, sometimes the agent resorts to pursuit of his or her self-interest instead of the interest that the principal had in mind. This is also known as the "agency" problem. According to PA theory, the agency problem arises when there is some sort of "information asymmetry" between the principal and the actor. In other words, the principal does not know what the agent is doing. The solution to overcome this asymmetry is obviously monitoring. However, principals cannot always monitor agents. The

commonest way of presenting this dilemma is represented by Grossman and Hart:

> Consider two individuals who operate in an uncertain environment and for whom risk sharing is desirable. Suppose that one of the individuals (known as the agent) is to take an action which the other individual (known as the principal) cannot observe. Assume that this action affects the total amount of consumption or money which is available to be divided between the two individuals. In general, the action which is optimal for the agent will depend on the extent of risk sharing between the principal and the agent. The question is: What is the optimal degree of risk sharing, given this dependence? (Grossman and Hart 1983: 7)

The analysis of risk-sharing is one approach in this debate. It has produced impressive models which focus on how to calculate costs and benefits – however at a quite abstract level. And of course it has produced a strong emphasis on how to handle the agency problem via contracts of various kinds (Holmström, 1979, Grossman and Hart 1983). Contracts and their design are a central feature of this discourse promoted to solve dilemmas of moral hazard, self-interest, and simply poor understanding of common objectives.[2]

There is certainly a claim here that PA theory offers important solutions by presenting a challenge to, for example, a public administration like this. It puts an emphasis on negotiations for contracts which will then steer processes forward with costs reduced as a consequence. Central problems associated with information asymmetries are tackled by setting up arrangements which may provide "credible commitments." (Groenendijk 1997; Geraldi 2007; Teorell 2007; Rothstein 2011b). Today PA theory holds a central position in administrative studies. It is generally portrayed as one of the pillars of what has commonly been known for a while as New Institutional Economics and it is regarded as one of the dominant theories in helping us explain why certain governments perform poorly and why they may be plagued by corruption (Rothstein 2011a; Persson et al., 2013).

## Objections to the Principal-Agent theory

PA theory is however by no means universally accepted or hailed. For example Persson, Rothstein and Teorell show a healthy scepticism towards PA theory and claim that in reality there are some places where corruption is so widespread that it may be hard to find the "principled principal", and then Collective Action perspectives on corruption are more useful (Rothstein 2011a; Persson et

2    For a comprehensive overview, see Laffont and Martimort 2002.

al., 2013). Nonetheless, PA theory is still a perspective which researchers tend to want to improve upon when deficiencies are brought up (Lane 2005; Persson and Sjöstedt 2012; Dehousse 2008). It is even the case that some public administration researchers regard PA theory as a substitute for NPM (Boston and Halligan 2012). Therefore, although a complete case will not be made here, let us at least take the critical analysis a few steps further by trying to formulate the main reasons why, to begin with, PA theory is not a substitute for NPM, and why it may be seen as overhyped theory in general.

The main objections raised here are that:
- PA theory has a reductionist perspective on incentives.
- The concepts of principal and agent are too abstract.
- PA theory works better as a metaphor than as a theory.

## PA theory, assumptions, incentives, and consequences

A substantial portion of the criticism made of PA theory is the same as the most common critical comments generally raised against Rational Choice theory. This is not a surprise since PA theories have roots in PC and RC theories. Consequently several objections that have been raised against RC theories apply to PA theory as well:

> The rational-choice model that has been subjected to the most criticism is one that presumes (1) extreme selfishness, (2) complete information, (3) an unambiguous capacity to attach utility to all outcome and actions and to rank all alternatives in a consistent manner, and (4) maximization of expected utility. (Rabin et al., 2007: 1102)

Here we recognize many core features of the criticism directed at the NPM model. The NPM model assumed that public servants were mainly motivated by selfishness, that they therefore needed to be monitored, and that evaluations could give complete pictures of activities and goal attainment which would therefore be the main instrument for steering organizations and individuals in the direction which would provide the maximum level of efficiency. All these assumptions have been important in the NPM models that have been applied since the late 1980s which gave rise to the evaluation explosion seen in the 1990s and which still lay a heavy burden on the public sector today (Power 1997; Brante et al., 2015).

It is also responsible for strategies based on monitoring which most employees regard as intrusive and having a de-professionalizing effect. What is puzzling here is that PA theory is based on many of the same ideas, but has man-

aged to survive far better in the academic literature than the widely criticized New Public Management models.

The point is not simply that PA theory is bad just because it builds on assumptions from Rational Choice theory. Selfishness can certainly be a strong motivating force. The perverse effects appear when RC postulates are applied in a one-sided manner. PA theory may appear to be a more sophisticated attempt as it points towards "information asymmetries" which can be more inclusive than just incorporation of "self-interest", but it still assumes that selfishness is the primary force, sometimes even the only force, which moves individuals. And then the main policy remedy which follows is monitoring and contracts. The way it does so is often quite deceptive, and it finally affects policy makers on the ground.

Based on what I have seen in this field of research and policy development, it happens like this: first, economists create elaborate models based on assumptions – in this case, for example, models that try to handle selfishness, information asymmetries, moral hazard, insurances etc. Models are then constructed that may provide illustrations that tell us for example that certain types of contracts may be needed after risk or asymmetry has reached a certain level. The level in question may not be possible to translate into something concrete in reality. However, an argument can still be provided about something we need to be aware of. Despite this, it should be remembered that the model which was the point of departure to illustrate interaction or interdependencies at a more abstract level was built on assumptions. Even if the economist making the model has a fair understanding of the fact that reality may be more complicated than the model assumes, the conclusions from the model can "catch on" and be elevated to a higher status in terms of claims on reality by those who want to present an appealing theory or policy. In this case it is PA theorists or NPM policy proponents. The theorists take the assumption made by the economists, and then regard it as a reality. By making this leap, PA theorists will then proceed to deny the importance of, for example, institutionalism as well as studies on mores and values such as those provided in the genre that Kaufman belongs to (Kaufman 2006). In the case of NPM policy makers, they project a view of employees as driven solely by self-interest, to a point where professionalism, honour, expertise, ethics, and other values are denied any space. Hence the emphasis put on monitoring, contracts and constant evaluations by defenders of both NPM and PA theory.

Even if this is a valid account of how assumptions travel to become facts and even policy proposals, it still happens even when the economists themselves warn against taking such leaps. Grossman and Hart, for example, clearly state that their models take into account only a situation where the principal cannot monitor the agent, and that that situation is far more complicated when "the agent possesses information about his environment /.../ which the principal

does not." (Grossman and Hart 1983). However, when a simplified theoretical concept captures the imagination of political scientists, economists, politicians, policy makers and public administrators then it is hard to make them let go.

The point here is not to claim that the whole field of public administration studies has become hypnotized by reductionist ideas provided by economists, or that all economists produce reductionist studies. Economists today wrestle more than ever with complex realities in, for example, the field of development economics. Economists like Sen, de Soto, Chattopadyay and Duflo are studying the basic realities of poverty and corruption in a more detailed way than many anthropologists (Sen 1999; de Soto 2000; Banerjee and Duflo 2011). A substantial part of the field of administration studies is not only freeing itself from the NPM paradigm but also doing more than resorting to elaborating on categories and concepts. For example, research in the field of crisis management studies shows the need for including perspectives provided by studies of social capital and political culture in order to understand what makes a public administration in particular, and crisis management institutions more generally, "do the right thing" (Benedict Dyson and 't Hart 2013; Persson et al., 2015; Boin et al., 2006). More thought is given to the fact that individuals are not only structural dopes who conform to the various institutional contexts they are exposed to. Causing change in the other direction, individuals can impose new norms and change organizational culture relatively quickly in a segment of the public administration.[3] Most important, in this genre of studies research now departs from the realization that administrative culture or shared norms can guide individual behaviour more than any set of formal rules, protocols, and white papers. This is a stark contrast to what the research field has looked like, and to what is the main guiding perspective of strategists in crisis management organizations in some countries today.

The claim here is that the problems of theorists and policy makers moving from assumptions to policy recommendations seem to be most pressing when we consider areas where political stakes are high – for example corruption studies, welfare politics, labour politics, and of course public administration policies. These areas provide impressive gravitational fields where once an issue is captured, fact, assumptions, normative claims and explanations can become quite difficult to disentangle.[4]

In sum, the general verdict in this section is that PA theory starts out from over-simplified ideas about what motivates behaviour, ideas which when applied in practice have already proved to generate perverse results (Hood and Peters, 2004). However, the simplistic rational choice assumptions that became

---

3    These considerations are the guiding ideas of the Persona project which focuses on crisis management in the EU. (URL: http://persona.statsvet.uu.se)

4    The inspiration for this description obviously comes from descriptions of *gravitational singularity*. So, we could call the situation described here *perceptional singularity*.

the pillar for NPM, and which are reproduced in PA theory, have pervaded the minds of a substantial number of policy makers, so it seems to have taken the role of a supra-ideology today. That is very hard to change.

The two remaining objections will be dealt with rather more briefly and to a large extent they follow from the criticism set out above.

## Abstractions and reality

As mentioned earlier, Persson, Rothstein and Teorell argue that in some contexts it can be hard to find anyone who actually fulfils the role of the "principled principal" that can be seen to represent "the public interest". It is mainly, we get from the description of Uganda and Kenya, a situation with a plethora of agents pursuing their own interests. Consequently, the authors argue, it is more useful to apply a collective action perspective to the corruption problems these two countries are facing. The position taken by Persson, Rothstein and Teorell is compelling, but the objections to PA theory raised are not necessarily detrimental to PA theory – unless spelled out more clearly.

PA theory simply states that "public interest" may be one interest that the principal can be said to represent. If we go to a country like Uganda we see that perspectives on interest vary. There are always situations where agency is a problem – no matter what interests are represented. It may be the military or the dominant party there: the National Resistance Movement. Its leaders all need to rely on someone else to do something for them. This is a part of the dilemma of delegating power which has been around as long as humans have interacted in more complex societies. So, to say that there are no principals championing "public interests" is not to say that Uganda's society is not full of agency problems – or principals. This is recognized by Persson, Rothstein and Teorell, but then the criticism is restricted to saying that PA theory may not be useful when studying "principled principals" in Uganda, while collective action theory may be so. Further problematization of the origins of PA theory is important here. Why did PA theory turn out to be a cul-de-sac when they took it to Uganda?

The more specific conclusion that there are few or no principals that represent a "public interest" is complicated. What does this conclusion tell us? Certainly Uganda's President Museveni would disagree with Persson, Teorell, and Rothstein. He sees himself as the father of the nation who saved it from the terror of Obote and Amin. That is the primary public interest that Museveni represents – at least officially. Also, consider General Katumba Wamala who always needs more funds to protect public interest by fighting "terrorists". And then we have Stanley Ntagali who is the Archbishop of the Anglican Church of in Uganda , who is advocating strict laws against homosexuals – in the name of the public interest. The important point here is that the concept of agency is blind to what the principal may represent – no matter which context we discuss. Uganda is just an exam-

ple. Naturally there are principals in Uganda. However, when we assume that a corruption-free administration is the main public interest then of course it may be hard to find leaders who champion this cause wholeheartedly in a country like Uganda – at least in the public administration, which to a large extent equals the main ruling party in the country: the National Resistance Movement – which is controlled by Museveni. However, the same leaders certainly see themselves as representing *other* public interests. Many of them would certainly appreciate it if politics and public services did not depend on bribes: but when weighing some public interests against others, corruption may be seen as a secondary problem. So, Persson, Teorell, and Rothstein may be right when they say that there may be a lack of principals defending the public interest of keeping a clean government and it may also be fruitful therefore to see corruption as a collective action problem there – which is a widely discussed point today in relation to the corruption debate in more general terms (Pippidi et al., 2011; Persson et al., 2013; Stephenson 2015; Marquette and Peiffer 2015). But what does the position taken by Persson, Teorell, and Rothstein really tell us about PA theory – other than that it does not seem to help when analysing Uganda? It seems that as soon as we move from the abstract models, provided mainly by the economists, and really apply the concept on the ground, close to empirical studies based on, for example, in-depth interviews, the complex realities make the assumptions brought in with the concept more of a burden than an asset. Dividing people into principals and agents is metaphorically compelling but, as Burden (Burden, 2007, 45) suggests, on closer scrutiny the dichotomy appears false.

Surely we can talk of agency in abstract terms, but if we then expect from it that one-dimensional characters will suddenly appear on the empirical map, allowing us to categorize actors as "principled principals" who always fight corruption, or selfish agents who always promote it, then it may be time for induction to replace deduction. What can we really expect to find if we imagine that the principal with only one true public interest objective in mind really exists? Or, as stated above, a world of agents driven mainly by self-interest? The problem here is that the concepts of principal and agent can be taken too far. Lane, and even Rose-Ackerman, point out that in reality it is challenging actually to decide or delimit who is the principal and who is the agent (Lane 2005). The corruption literature often emphasizes the fact that it is not only a matter of formally deciding where the boundaries go. In reality the roles of the principal and agent simply shift too rapidly. People are also connected to each other in so many ways that we cannot even expect to disentangle such relationships in a meaningful manner – at least not if we assume the role labels reflect reality. So, while the dichotomy between principal and agent is convenient from a modelling perspective – as we discussed earlier – it may actually make far too many simplifications about what actors represent and how isolated they are in relation to other interests. Whether this really is a problem for PA theory is

most importantly connected to whether we should call PA theory a theory or not. This takes us to the third and final point made here.

## The theoretical value of PA theory

It is certainly accepted here that using the term "agency" or "Principal–Actor problem" is sometimes quite appropriate in order to simplify analytical discussions so they can more quickly cut to the chase. The objection which is raised here in this final section on criticism is that it may be inappropriate to speak of a PA *theory* unless the proponents of PA theory manage to come up with something more tangible than pointing towards "selfishness" and "information asymmetry" in rather broad terms.

As stated above, it is important to weigh in self-interest when discussing what makes people do the right thing, but it is a narrow starting point which is sometimes taken so far that it provides only circular evidence. Public administration researchers, especially those with a comparative perspective, try hard to understand motivation. Motivation may come from a spectrum of sources and rewards. The factors range from those at the personal level to others at the contextual level. They also range from those that emphasize pure self-interest to altruistic motives relating to the will to serve the public interest, professional pride or even nationalistic goals. The public choice advocates however tend to interject here that the goals that we are apt to describe as the most altruistic are actually driven solely by self-interest and self-satisfaction (Perry and Recascino Wise 1990; Chowdhury 2011). However, a theory which claims to explain everything violates the basic condition of a theory – that it must be falsifiable. A theory which explains everything in a world of variation explains nothing.[5]

Of course many PA theorists may not be so one-sided that they do not realise this. But if we then lean more on the other PA leg – the claim about information asymmetry – hoping to find something more valuable here, we may once again end up disappointed. The discussion about information asymmetry is mainly descriptive. It essentially points out that the principal does not know what the agent does all the time and it discusses some forces relating to self-interest that may pull the agent away from the one of the principals. The asymmetry issue is the closest that PA theory comes to pointing to a mechanism at work which allows for corruption, and the solution is most often framed in terms of contracts of various designs. However, the emphasis on information asymmetry is quite general and has been pointed out as far back as we have written records. Most explicitly it was analysed by Kautiliya, the master bureaucrat working for the Gupta Empire, in his writings in the Arthaśāstra around 300 BCE:

---

5     Some such basic criteria for a theory are well explained for example in Popper 1963.

> Just as fish moving inside water cannot be known when drinking water, even so officers appointed for carrying out works cannot be known when appropriating money.

Kautiliya also stated that:

> It is possible to know even the path of a bird flying in the sky, but not the ways of officers moving with their intentions concealed.

So it is hard to make the case that the component regarding "information asymmetry" is one of the great discoveries of PA theorists. However what is more troubling is the lack of theoretical clarity from PA theorists who so strongly emphasize principals, agents, information asymmetry and contracts.[6] If this description is not too ungenerous, it is not clear what the PA literature actually provides which warrants calling it a theory.

## Epilogue

There seems to be some sort of consensus among researchers in administration studies that New Public Management models and theories have not only seen their best days, but they have also caused much damage in public administration and continue to do so. We are now at a stage where expectations arise that public administration researchers should fill the vacuum after NPM has lost its popularity even among political actors. However, it is important to be very careful before we accept new "new" theories in this area such as "value based" models or proponents of "New Political Governance" (Bakvis and Jarvis 2012). We need to remember that one of the most important lessons from the failures of NPM stemmed from the expectation that it could fit any problem encountered by the whole range of different organisations that a modern system of public administration consists of today. We have also seen that in spite of all this harsh criticism of NPM, even some of its own critics seem to be prone to gravitate towards the same old solutions in different clothing – namely the Principal Agent perspective. PA theory encourages what NPM did to create dysfunctions: it promotes the idea that most problems can be solved by regarding humans as mainly motivated by self-interest, by generating more information (to counteract the asymmetries) and detailed contracts, and it supports top-down control instead of promoting trust, personal capabilities, and professionalism. Nonetheless, the basic principles that shaped NPM still exert such power over the mind that it is hard for many to think outside its limited universe. Even historians are so captivated by some of the related ideas that this is the only lens

---

6    See for example Sjöstedt (2009) who provides a very interesting explanation regarding "credible commitments" relating to this discussion and which does so well without PA theory.

which is applied not only when we look ahead towards what needs to be done, but also when looking backwards trying to understand what happened.[7]

So it has been advocated here that if we now also take PA theory out of the back-pack, moving forward will be easier. Policy makers will find it easier to adapt to the different terrains of the increasingly complex landscapes they encounter, researchers can broaden the empirical analysis, and, hopefully, fewer employees in the public sectors will get hurt. It is important to realize that the public sector cannot be managed by reference to the economist perspective alone. In saying so, the point is not to drop RC theory. RC theory can be successfully applied the way it is done by, for example, Collective Action theorists – in combination with other theories – especially those which emphasize social capital and political and administrative culture studies (without of course being overtaken by another reductionist idea saying that "everything is culture").

The main point which is made here is that public administration studies in some forms, public administration policy makers and political actors, have allowed themselves to accept perspectives on complicated problems and challenges in public administration which may appear to be qualified, sophisticated and technically advanced, whereas upon closer acquaintance they prove to be quite reductionist and to build on over-simplifications that are highly misleading. It is a systematic interpretation of complex problems as reducible to fairly simple incentives. This is not only sad. It implies a dangerous submission to powers which invert the role of science and theory. The ideal of the social *scientific* enterprise is to reveal the true nature of a phenomenon; first to understand its causes, before we even begin to speak in terms of policies, remedies and recommendations. However the reductionism we have witnessed in New Public Management, Public Choice, and ultimately in the Principal Agent perspectives discussed here, and how these perspectives have been applied, do the opposite. Most likely the greatest charge against these perspectives and movements is that they have made far-reaching and simplistic assumptions about human behaviour, and then not only sold it as theory. They have gone as far as actually imposing their views on human behaviour so strongly that human behaviour has been forced to conform to the stipulated parameters. This is what happened when NPM-based theories advocated constant and simplistic evaluations. The theoretical assumptions which steered the evaluations made humans behave accordingly. For example, if the evaluations decided to measure quality in the care of the elderly solely by setting twelve or fourteen minutes for a shower every second or third day as a "norm" for each "customer"– then this would establish a practice where no elderly person would ever get more than that in terms of support for personal hygiene. Having a cup of coffee and a sim-

---

7    This is a common criticism of, for example, Francis Fukuyama and the perspectives projected in "The end of history and the last man" (1992). See for example Glaser (2014) and Stanley and Lee (2014).

ple chat about the weather would certainly fall outside the scope of efficiency and measurable qualities (Zaremba, 2013). When the only incentive offered as a reward for public servants was expressed in terms of salary, everything but salary would be pushed out – including professionalism. To put it bluntly, economic theories, and their adherents, are far too little concerned with taking broad perspectives on how humans think, feel and react. They are far too much concerned with re-modelling human behaviour, in order to bring it into line with their own theoretical image of it. This is the paradigm that needs to be broken in order to move forward in a post-NPM world.

# References

Ahlbäck Öberg, S., 2010. "Att kontrollera förvaltningen: framväxten av granskningssamhället", in Rothstein, Bo (ed.), *Politik som organisation: förvaltningspolitikens grundproblem* (pp. 168–196). Stockholm: SNS Förlag.

Ahlbäck Öberg, S. & Bringselius, L., 2015. "Professionalism and organizational performance in the wake of new managerialism", *European Political Science Review,* 7(4), 499–523.

Ahlbäck Öberg, S. & Widmalm, S., 2012. "Professionalism nedvärderas i den marknadsstyrda staten", *Dagens Nyheter*, 2012-10-26.

Ahlbäck Öberg, S. & Widmalm, S., 2013. "NPM på svenska", in Zaremba, M. (ed.), *Patientens pris.* Stockholm: Weyler förlag.

Bakvis, H. & Jarvis, M. D. (eds), 2012. *From New Public Management to New Political Governance.* McGill-Queen's University Press.

Banerjee, A. V. & Duflo, E., 2011. *Poor Economics – A Radical Rethiniking of the Way to Fight Global Poverty.* New York: Public Affairs.

Benedict Dyson, S. & 'T Hart, P., 2013. "Crisis Management", in Huddy, L. S., David O., Levy, Jack S. (ed.), *The Oxford Handbook of Political Psychology.* Oxford University Press.

Boin, A., 'T Hart, P. & Sundelius, B., 2006. *The Politics of Crisis Management – Public Leadership Under Pressure.* Cambridge: Cambridge University Press.

Boston, J. & Halligan, J., 2012. "Political Management and New Political Governance: Reconciling Political Responsiveness and Neutral Competence", in Bakvis, H. & Jarvis, M. D. (eds.), *From New Public Management to New Political Governance.* Montreal & Kingston: McGill-Queen's University Press.

Brante, T., Johansson, E., Olofsson, G. & Svensson, L. G., 2015. *Professionerna i kunskapssamhället.* Stockholm: Liber.

Buchanan , J. M. & Tullock, G., 1962. *The Calculus of Consent: Logical Foundations of Constitutional Democracy.* Minneapolis: Liberty Fund.

Burden, B., 2007. *Personal Roots of Representation.* Princeton: Princeton University Press.

Chowdhury, N., 2011. "Principal- Agent Theory Is an effective Tool or only a good Hypothesis?", *Development Management*, University of Birmingham.

De Soto, H., 2000. *The Mystery of Capital.* London: Black Swan.

Dehousse, R., 2008. "Delegation of powers in the European union: The need for a multi-principals model", *West European Politics,* 31, 789–805.

Dunleavy, P. & Hood, C. 1994. "From Old Public-Administration to New Public Management", *Public Money & Management,* 14, 9–16.

Fukuyama, F., 1992. *The end of history and the last man.* London: Hamish Hamilton.

Geraldi, J. G., 2007. "New Institutional Economics". *Fachbereich Maschinenbau – Management internationaler Projekte.* Universität Siegen.

Glaser, E., 2014. "Bring back ideology: Fukuyama's 'end of history' 25 years on", *The Guardian*, 21 March.

Groenendijk, N., 1997. "A principal-agent model of corruption", *Crime, Law & Social Change,* 27, 207–229.

Grossman , S. J. & Hart, O. D., 1983. "An Analysis of the Principal-Agent Problem", *Econometrica* 51, 7–45.

Holmström, B., 1979. "Moral Hazard and Observability", *The Bell Journal of Economics,* 10, 74–91.

Hood, C., & Peters, G., 2004. "The Middle Aging of New Public Management: Into the Age of Paradox?", *Journal of Public Administration Research and Theory: J-PART*, *14*(3), 267–282

Jensen , M. C. & Meckling, W. H., 1976. "Theory of the Firm: Managerial Behavior, Agency Costs and Ownership Structure", *Journal of Financial Economics,* October, 1976, Vol. 3, No. 4, 3, 305–360.

Kaufman, H., 2006. *The Forest Ranger – a study in administrative behavior.* New York: Resources for the Future.

Klitgaard, R., 1988. *Controlling Corruption.* Berkeley: University of California Press.

Kuhn, T., S., 1970. *The Structure of Scientific Revolutions.* Chicago: The University of Chicago Press.

Laffont, J.-J. & Martimort, D., 2002. *The Theory of Incentives: The Principal-Agent Mode.* Princeton: Princeton University Press.

Lane, J.-E., 2005. *Public Administration & Public Management: The Principal-Agent Perspective.* London: Routledge.

Liedman, S.-E., 2012. *Hets!: – En bok om skolan.* Stockholm: Albert Bonniers Förlag.

Lindgren, L., 2014. *Nya utvärderingsmonstret – Om kvalitetsmätning i den offentliga sektorn.* Lund: Studentlitteratur.

Marquette, H. & Peiffer, C., 2015. "Corruption and Collective Action". *Development Leadership Programme, and Anti-Corruption Resource Centre.*

Perry, J., L . & Recascino Wise, L., 1990. "The Motivational Bases of Public Service", *Public Administration Review,* 50, 367–373.

Persson, A., Rothstein, B. & Teorell, J., 2013. "Why Anticorruption Reforms Fail– Systemic Corruption as a Collective Action Problem", *Governance: An International Journal of Policy, Administration, and Institutions,* 26, 449–471.

Persson, A. & Sjöstedt, M., 2012. "Responsive and Responsible Leaders: A Matter of Political Will?", *Perspectives on Politics,* 103, 617–632.

Persson, T., Widmalm, S. & Parker, C., 2015. *Social trust, impartial administration and public confidence in EU crisis management institutions.*

Pippidi, A. M., Loncaric, M., Mundo, B. V., Braga, A. C. S., Weinhardt, M., Solares, A. P., Skardziute, A., Martini, M., Agbele, F., Jensen, M. F., Soest, C. V. & Gabedava, M., 2011. "Contextual Choices in Fighting Corruption", in Pippidi, A. M. (ed.)., Berlin: Norad, c/o ANKOR (the Anti-corruption Project), in cooperation with the Evaluation Department ("Contextual Choices for Results in Fighting Corruption").

Pollitt, C., 2014. "Managerialism redux? Keynote address to the 2014 EIASM conference, Edinburgh", 2014-08-24 [https://soc.kuleuven.be/io/nieuws/managerialism-redux.pdf, åtkomst 2015-10-21].

Popper, K. R., 1963. *Conjectures and Refutations.* London: Routledge.

Power, M., 1997. *The Audit Society – Rituals of Verification.* Oxford: Oxford University Press.

Rabin, J., Hildreth, B. W. & Miller, G. J., 2007. *Handbook of Public Administration.* Boca Raton, CRC Press Taylor and Francis Group.

Rose-Ackerman, S., 1999. *Corruption and Government.* Cambridge: Cambridge University Press.

Rothstein, B., 2011a. "Anti-corruption: the indirect 'big bang' approach", *Review of International Political Economy,* 18, 228–250.

Rothstein, B., 2011b. *The Quality of Government.* Chicago: Chicago University Press.

Sen, A, 1999. *Development as Freedom.* Oxford: Oxford University Press.

Sjöstedt, M., 2009. *Thirsting for Credible Commitments How Secure Land Tenure Affects Access to Drinking Water in Sub-Saharan Africa.* PhD, Göteborg University.

Stanley, T. & Lee, A., 2014. "It's Still Not the End of History", *The Atlantic.*

Stephenson, M., 2015. "Corruption is BOTH a 'Principal-Agent Problem' AND a 'Collective Action Problem'". Available from: http://globalanticorruptionblog.com/2015/04/09/corruption-is-both-a-principal-agent-problem-and-a-collective-action-problem/ [Accessed April 9 2015].

Teorell, J., 2007. "Corruption as an Institution Rethinking the Nature and Origin of the Grabbing Hand", *QoG Working Paper Series.* Gothenburg Quality of Governance, Gothenburg University.

Tingsten, H., (1945) 1960. *Demokratiens Problem.* Stockholm: Aldus, Bonniers.

Widmalm, F., Öberg, S. A. & Widmalm, S., 2013. "Vårdens kontrollsystem bättre för andra områden", *Dagens Nyheter,* 2013-04-17.

Widmalm, S., 2012. "Utvärdering blir ytvärdering", *Respons,* 11–13.

Widmalm, S. & Ahlbäck Öberg, S., 2013. "NPM på svenska", in Zaremba, M. (ed.), *Patientens pris – Ett reportage om den svenska sjukvården och marknaden.* Stockholm: Svante Weyler Bokförlag AB.

Widmalm, S., Widmalm, F. & Persson, T., 2014. "Välfärden undermineras av vårt fokus på utvärderingar", *Dagens Nyheter,* 2014-04-19.

Zaremba, M., 2013. *Patientens pris – Ett reportage om den svenska sjukvården och marknaden.* Stockholm: Svante Weyler Bokförlag AB.

# THE INFORMATION ASYMMETRY ASPECT OF AGENCY THEORY IN BUSINESS COMPLIANCE CONTEXTS: A SYSTEMATIC REVIEW

**Omar Omar[1], Denilson Sell[2], Aires José Rover[3]**

**Abstract.** *The information asymmetry aspect of agency theory constitutes a relevant risk and hinders principal-agent relationships. Business compliance implementations, therefore, must address the information asymmetry aspect. This article summarizes the existing literature concerning information asymmetry, agency theory and compliance through a systematic literature review in the Web of Science and Scopus databases. A broad set of applicable domains was identified, with varied approaches to tackle information asymmetry and mixed results. Firms with successful compliance implementations, and those who implement compliance beyond obligatory regulation achieve higher levels of performance. The systematic review also indicates promising research opportunities addressing the convergence of knowledge engineering artefacts, possibilities and knowledge management tools towards business compliance and the information asymmetry of the principal-agent problem.*

**Keywords:** *information asymmetry; agency theory; principal-agent problem; compliance; knowledge engineering*

---

[1] Graduate Program of Knowledge and Engineering Management – Federal University of Santa Catarina (UFSC) Florianópolis – SC – Brazil. Email: omarx02@gmail.com
[2] Professor at the Graduate Program of Knowledge and Engineering Management – Federal University of Santa Catarina (UFSC) Florianópolis – SC – Brazil. Email: sell@stela.org.br
[3] Professor at the Graduate Program of Knowledge and Engineering Management – Federal University of Santa Catarina (UFSC) Florianópolis – SC – Brazil. Email: aires.rover@gmail.com

## 1    INTRODUCTION

This article is part of broader research concerning agency theory in corporate compliance contexts. Our research hypothesis is that knowledge engineering methods and tools can help mitigate the principal-agent problem and address the information asymmetry gap between the principal and agent. The research's final goal is to propose a knowledge model to be used in information intensive business compliance contexts. We analyze business compliance from the principal-agent perspective, in which the principal-agent relationship is characterized by information asymmetry, uncertainty and opacity (Taleb, 2012 apud Omar et al, 2016). A thorough study of the information asymmetry aspect of agency theory and compliance is, therefore, required; and it is the object of this study.

Business compliance can be seen as a set of rules and codes implemented to ensure the fulfillment of business objectives in a transparent and correct manner; mitigating risk and ensuring that the organization complies with respective, both internal and external, regulation (Kuhnel et al, 2017; Marekfia et al, 2012; Pupke, 2008). Within the enterprise, the relationship between principals (mainly stakeholders and hierarchical decision makers) and agents (executives and staff) constitutes one of business compliance's core and main concerns (Kuhnel et al, 2017; Pupke, 2008). Thus, agency theory constitutes an adequate and useful tool to interpret and analyze business compliance.

The information asymmetry aspect of agency theory jeopardizes the principal-agent relationship, which present an inherently principal-agent problem (Hoenen; Kotsova, 2015; Eisenhardt, 1989). Agency theory is analog to information processing approaches of contingency theory, as the two perspectives constitute information theories (Silva, 2016; Eisenhardt, 1989). Within this context, considering that both principal and agent are rational and seek to maximize their utility functions, often having nonaligned and different interests, and that they possess different pieces of, and access to, information, the principal-agent problem can be understood as an information asymmetry and uncertainty problem (Silva, 2016; Hoenen; Kotsova, 2015; Saito; Silveira, 2008; Eisenhardt, 1989).

Compliance is crucial to ensure that the principal interests are protected, carried on and cared for by the agents. Understanding the relationship between compliance and the information asymmetry aspect of agency theory through the lens of existing scientific literature is therefore paramount to better address business compliance initiatives and planned or ongoing implementations, contributing to the satisfactory fulfillment of compliance requirements.

## 2    METHODOLOGY

Systematic literature review is a method that allows the summarization of large quantities of contents of a given subject (Cook; Murlow; Haynes, 1997) and was the method employed in this article. According to Cordeiro et. al (2007), systematic reviews provide credibility and robustness to collected data because they employ rigorous and explicit steps to identify, select, collect and analyze data, and to describe the relevant contributions of these data to a given research.

To Sampaio and Mancini (2007) apud Omar, Cunha and Sell (2016), systematic reviews must: a) be reproducible; b) present a crystal clear research question; c) employ defined search strategies; d) possess inclusion and exclusion criteria; e) carefully analyze found literature; and f) follow an original article structure.

The review process was conducted according to Rother (2007) apud Botelho, Cunha e Macedo (2011) proposed seven steps: a) formulation of the research question; b) tracking of the articles; c) critical evaluation of the articles; d) data collection; e) data analysis and presentation; f) data interpretation; and g) review update and enhancement.

The research question that guided the systematic review was: "what research has been conducted concerning the information asymmetry aspect of agency theory in business compliance contexts?". Other assumptions are that information asymmetry can hinder compliance in a wide range of domains and that knowledge tools can help mitigate its effect on compliance.

Initial search with the keywords "agency theory"; "information asymmetry"; "corporate compliance"; and/or "business compliance" has returned few results, indicating few existing literature regarding the research subjects. Additional searches adding the keywords "principal-agent problem" instead of "agency theory" and broader scope using the keyword "compliance" by itself were performed.

Once the keywords were determined, searches were performed using Web of Science and Scopus databases, employing logical operators to alternate and/or concatenate keywords. The criteria employed to include or not a publication were the abstract and article analysis and subsequent adherence to the research question. The timespan searched is between 2008 and 2017. Eleven articles matched these criteria and were selected; as listed in table 1.

Table 1 – selected articles for the systematic review ordered by # of citations

| Article | Journal | Author (s) | Year | Citations |
|---------|---------|------------|------|-----------|

| # | Title | Journal | Authors | Year | Citations |
|---|-------|---------|---------|------|-----------|
| 1. | Why do firms go dark? Causes and economic consequences of voluntary SEC deregistrations | Journal of Accounting and Economics | Leuz, Christian; Triantis, Alexander; Wang, Tracy Yue | 2008 | 86 |
| 2. | To what extent are EU steel companies susceptible to competitive loss due to climate policy? | Energy Policy | Okereke, C., McDaniels, D. | 2012 | 19 |
| 3. | Aggregated, voluntary, and mandatory risk disclosure incentives: Evidence from UK FTSE all-share companies | International Review of Financial Analysis | Elshandidy, Tamer; Fraser, Ian; Hussainey, Khaled | 2013 | 17 |
| 4. | The impact of internet health information on patient compliance: A research model and an empirical study | Journal of Medical Internet Research | Laugesen, J., Hassanein, K., Yuan, Y. | 2015 | 12 |
| 5. | Sustainability in multi-tier supply chains: Understanding the double agency role of the first-tier supplier | Journal of Operations Management | Wilhelm, M.M., Blome, C., Bhakoo, V., Paulraj, A. | 2016 | 10 |
| 6. | Corporate Legitimacy and Investment-Cash Flow Sensitivity | Journal of Business Ethics | Attig, N., Cleary, S.W., El Ghoul, S., Guedhami, O. | 2014 | 6 |
| 7. | The impact of external monitoring and public reporting on business performance in a global manufacturing industry | Business and Society | Katz, J.P., Higgins, E., Dickson, M., Eckman, M. | 2009 | 6 |
| 8. | Exploring agency, knowledge and power in an Australian bulk cereal supply chain: A case study | Supply Chain Management | Byrne, R., Power, D. | 2014 | 5 |
| 9. | Extending the boundaries of IQ: Can collaboration with information management improve corporate governance | International Journal of Information Quality | Maguire, H. | 2008 | 2 |
| 10. | Information security policy compliance: An empirical study on escalation of commitment | 19th Americas Conference on Information Systems, AMCIS 2013 - Hyper connected World: Anything, Anywhere, Anytime | Kajtazi, M., Bulgurcu, B. | 2013 | - |
| 11. | Agency Theory, Disclosure - Transparency: The Nemesis of Enterprise and Corporate Governance Systems | 4th European Conference on Management, Leadership and Governance (ECMLG) | Lazarides, Themistokles; Argyropoulou, Maria; Drimpetas, Evaggelos; | 2008 | - |

Source: compiled by the author.

The next section discusses the results of the articles' analysis.

# 3      ANALYSIS OF SELECTED ARTICLES

All the selected articles consider at least two elements of the research subjects of this research: the relationship between compliance (in a broader sense, not limited to business compliance) and information asymmetry and/or agency theory or the principal-agent problem.

One of the first conclusions is that the subject and its' relationship with the elements of the database search performed is little explored by the respective scientific communities. Several wide results were achieved when employing a single search keyword in either case (compliance or agency theory), but the fact is that few studies have addressed both as a combined approach to address the information asymmetry issue. Furthermore, few studies address information asymmetry as a relevant aspect towards the success or failure of compliance implementations.

Nonetheless, the results showed a broad range of domains affected by information asymmetry between principals and agents, strengthening our initial supposition.

Among the domains covered by the articles are US public trading firms that stop SEC reporting (Leuz et al, 2008); EU steel industry and incentives to sustain competitiveness while complying with environmental compliance (Okereke; Mcdaniels, 2012); public non-financial trading companies of the FTSE in the UK, and their risk disclosure policies (Elshandidy et al, 2013); medical domains and the public information on the internet regarding health and diseases (Laugesen et al, 2015); supply chains and sustainability compliance throughout the supply chain and the role of first tier suppliers as double agents of sustainability compliance (Wilhelm et al, 2016); Investment cash flow sensitivity and the impact of corporate social responsibility and legitimacy on investment cash flow (Attig et al, 2014); apparel industry as a global manufacturing industry and the impact of public reporting and external monitoring (Kats et al, 2009); bulk commodities (cereal) Australian industry supply chains and the role of knowledge to form collaborative supply chains (Byrne; Power; 2014); corporate governance and collaboration together with information management to enhance compliance (Maguire, 2008); information security policy and behavioral study of employees' compliance with escalation rules (Kajtazi; Bulgurcu, 2013); and a theoretical study of the relationship between agency theory, corporate governance and information systems (Lazarides et al, 2008).

Table 2, listed below, summarizes the selected articles' main topics and research domains and their respective results and impact on information asymmetry.

Table 2 – Main Topics and results

| Article # | Main Topics / research domain | Main results / Impact of information asymmetry |
|---|---|---|
| 1 | US public trading firms that cease SEC reporting (hence going dark); main causes are the SOX act (compliance); poor future prospects; and distress. Published soon after the adoption of the Sarbanes-Oxley Act. | Compliance increased costs cause the firms to go dark; going dark jeopardizes the firms' returns; Insiders have vested interests in going dark, decreasing outside scrutiny and increasing the information asymmetry, happens more often when governance and investor (principal) protection are weak; Going dark and going private are two different things. |
| 2 | Steel and Iron industries within the European Union; Specifically, ArcelorMittal; Corus; and ThyssenKrupp. Special rules devised to increase the sector's competitiveness regarding carbon footprint. Highlights the information asymmetry aspect of the relationship between the government (principal) and agents (steel firms) | The steel companies largely exaggerate their competitive vulnerabilities in order to receive greater incentives and free allocations/allowances. This largely happens due to the principal-agent problem (information asymmetry aspect specifically) of the companies' narratives regarding their vulnerabilities and the EU ETS inability to objectively probe and assess their competitive vulnerabilities. |
| 3 | Analyzes the impact of corporate risk levels on aggregates, voluntary and mandatory risk disclosures of FTSE non-financial listed companies; The article contrasts firms with greater compliance with mandatory regulations with those of lower compliance requirements. | The article draws a positive correlation between higher board independence and lower insider/outsider interference with overall firm returns and higher dividend yields; Higher levels of risks are related to higher levels of voluntary risk disclosures (reducing information asymmetry and increasing compliance); |
| 4 | A survey study of 225 medical patients that use health information on the internet and its impact on the patients' compliance with the prescribed treatment regimens; The article analyzes the impact of the information asymmetry (of using information on the internet) in the relationship of patient (agent) and doctor (principal). | Internet public health information has little to no impact on patient compliance with treatments; The principal (physician) has the upper hand due to the patient's perceptions of higher quality information; This article offers an uncommon perspective where information asymmetry acts to the benefit of the principal instead of the agent. |
| 5 | Dispersed multi-tier supply chains with great and growing complexity. The article states that first-tier suppliers are crucial for sustainability compliance throughout the supply chain. First-tier suppliers act as agents and double agents (forwarding requirements to their own supply chains) to fulfill the leading firm's (principal) sustainability compliance requirements. Three case studies of different institutional contexts were carried. | Major findings include the leading firm's own internal alignments of the sustainability function determine the supply chain compliance; reducing the information asymmetry along the supply chain positively contributes to the supply chain sustainability compliance; The authors propose a framework for future research regarding multi-tier supply chains; and also highlight the double agency role of first tier suppliers as an instrument to fulfill compliance requirements. |
| 6 | Assesses the impact of corporate social responsibility (CSR) on investment sensitivity to cash flows (ISCF). Information asymmetry and agency costs of CSR affect investment-cash flow sensitivity. | CSR performance decreases ISCF; ISCF increases when CSR concerns increase; CSR activities beyond compliance requirements is desired and may improve access to financial capital. |
| 7 | Focuses on external monitoring and public reporting and its correlation with firms performance in terms of valuation in the apparel industry domain. | External monitoring is valuable to business (agent) and society (principal) because it reduces the information asymmetry mitigating the principal-agent problem; |

| | | |
|---|---|---|
| 8 | Inter-firm relationships and the role of information sharing practices in bulk commodity supply chains in Australia, characterized by power asymmetries in the system. Contrasts compliance x collaboration of different firms in the supply chain | Better understanding of the relationship between firms translates in better procurement and better chances of creating situations of collaboration instead of plain compliance; |
| 9 | Agency theory moral hazards, adverse selection and consequently information asymmetry between CEOs and boards, coupled with widespread and overwhelming information volumes within the enterprise increase the pressure on the firm to meet legislative, accountability, and business requirements and compliance. | Effective records management and record management functionality to asses information quality can improve corporate governance mechanisms and meet the requirements of accountability, transparency and compliance and in reducing information asymmetry. |
| 10 | The authors work within the boundaries of agency theory and the theory of planned behavior to examine Escalation situations and employees' behavior towards compliance with the requirements of their information security policies. | The authors delineate three mediating factors to explain employees' attitude: work impediment; information asymmetry; and safety of resources; Main finding is that information asymmetry and safety of resources both have significant impact on attitude, but work impediment does not; information security policy, also, as significant impact on all mediating factors. |
| 11 | Discusses information requirements for compliance with SOX act and OECD's corporate governance principles. Information flow as a critical factor for the success of corporate governance and depends on the firm's information system's design. | Corporate Governance depends on the implementation of modern enterprise systems to secure disclosure and transparency; The control and dissemination of information are crucial to comply with OECD principles and SOX. |

Source: Sourced from the selected articles (table 1).

Table 2 analysis evidences that despite having different focus and approaches, all articles highlight the information asymmetry problem and the challenges posed by it when implementing compliance requirements.

Two articles cover the subject in the supply chain domains (Wilhelm et al, 2016) and (Katz et al, 2009). Supply chains constitute a relevant compliance domain due to the many distinct independent actors involved, and risks posed, whether from the point of view of supply interruption and business disruption or the dependencies on supply companies to comply with a wide range of requirements, from sustainability and environmental to materials used and even labor policy and financial transparency. The issue is especially relevant in multi-tier supply chains (Wilhelm et al, 2016). Understanding the agency role of suppliers and addressing the principal-agent problem both constitute determinant factors to successfully achieve compliance throughout the supply chain.

Other articles focus on public listed companies and information disclosure, whether it is mandatory or voluntary disclosure, and even the absence of it - firms going dark – as covered by Leuz et al (2008) in the US financial market; or risk disclosure incentives from public traded FTSE companies in the UK (Okereke; McDaniels, 2012). In both cases, the authors address the principal-agent problem from the perspective of the stock market and society in general (principal) and the

736

companies per se (agent) and highlight the constraints and limitations imposed by information asymmetry.

Several other articles discuss the information asymmetry problem through the realities of several distinct domains. And although most articles propose some advice or recommendations on how to address and mitigate the issue of information asymmetry, none have resorted to knowledge management tools or knowledge engineering artefacts to do so. Moreover, a lack of a structured approach, possibly due to the uncertain and implicit nature of principal-agent relationships, to address the information asymmetry issue was also observed.

## 4    FINAL CONSIDERATIONS

One of the results of the systematic review is the fact that there are few articles in the existing literature that correlate agency theory, information asymmetry and compliance, indicating a promising field of original and relevant research. This observation is especially apparent regarding business compliance, specifically.

While analyzing the selected articles of the systematic review, the lack of knowledge engineering or knowledge management approaches to tackle the information asymmetry aspect of the principal-agent problem became evident. We believe that due to the nature of the information asymmetry problem, information systems and knowledge methods, artefacts and tools would yield positive outcomes tackling the issue, and constitute a promising research field.

The implicit nature of the asymmetry of information within the principal-agent relationship and its inherent complexity poses a challenge to the elicitation, structuring and sharing of information and knowledge throughout the enterprise and towards external stakeholders, and we believe that knowledge engineering has a relevant role to play in this context.

It is important to emphasize that the articles that somehow showed cases in which the information asymmetry problem was addressed and mitigated regardless of method were able to achieve higher performance and overall better competitiveness, indicating that the enterprise benefits from compliance, which represent a substantial additional incentive for implementing it. Nonetheless, objective measurements of these benefits and quantitative assessment of the results of compliance, even beyond imposed regulations, are also promising fields of research.

## 5    REFERENCES

Attig, N., Cleary, S. W., El Ghoul, S., & Guedhami, O. (2014). Corporate legitimacy and investment–cash flow sensitivity. Journal of Business Ethics, 121(2), 297-314.

Botelho, L. L. R., Cunha, C. C. D. A., & Macedo, M. (2011). O método da revisão integrativa nos estudos organizacionais. Gestão e Sociedade, 5(11), 121-36.

Byrne, R., & Power, D. (2014). Exploring agency, knowledge and power in an Australian bulk cereal supply chain: a case study. Supply Chain Management: An International Journal, 19(4), 431-444.

Cook, D. J., Mulrow, C. D., & Haynes, R. B. (1997). Systematic reviews: synthesis of best evidence for clinical decisions. Annals of internal medicine,126(5), 376-380.

Cordeiro, A. M., Oliveira, G. M. D., Rentería, J. M., & Guimarães, C. A. (2007). Revisão sistemática: uma revisão narrativa. Rev. Col. Bras. Cir,34(6), 428-431

Eisenhardt, Kathleen M. "Agency theory: An assessment and review." Academy of management review 14.1 (1989): 57-74.

Elshandidy, T., Fraser, I., & Hussainey, K. (2013). Aggregated, voluntary, and mandatory risk disclosure incentives: Evidence from UK FTSE all-share companies. International Review of Financial Analysis, 30, 320-333.

Hoenen, A. K., & Kostova, T. (2015). Utilizing the broader agency perspective for studying headquarters–subsidiary relations in multinational companies. Journal of International Business Studies, 46(1), 104-113.

Kajtazi, M., & Bulgurcu, B. (2013). Information Security Policy Compliance: An Empirical Study on Escalation of Commitment.

Katz, J. P., Higgins, E., Dickson, M., & Eckman, M. (2009). The impact of external monitoring and public reporting on business performance in a global manufacturing industry. Business & Society, 48(4), 489-510.

Koufopoulos, D. (2008). Agency Theory, Disclosure-Transparency: The Nemesis of Enterprise and Corporate Governance Systems. In ECMLG2008-Proceedings of the 4th European Conference on Management Leadership and Governance: ECMLG (p. 103). Academic Conferences Limited.

Kühnel, S., Sackmann, S., & Seyffarth, T. (2017). Effizienzorientiertes Risikomanagement für Business Process Compliance. HMD Praxis der Wirtschaftsinformatik, 1-22.

Laugesen, J., Hassanein, K., & Yuan, Y. (2015). The impact of internet health information on patient compliance: a research model and an empirical study. Journal of medical Internet research, 17(6).

Leuz, C., Triantis, A., & Wang, T. Y. (2008). Why do firms go dark? Causes and economic consequences of voluntary SEC deregistrations. Journal of Accounting and Economics, 45(2), 181-208.

Lopes, A. L. M., & Fracolli, L. A. (2008). Revisão sistemática de literatura e metassíntese qualitativa: considerações sobre sua aplicação na pesquisa em enfermagem. Texto & Contexto-Enfermagem, 17(4), 771-778.

Maguire, H. (2008). Extending the boundaries of IQ: Can collaboration with information management improve corporate governance. International Journal of Information Quality, 2(1), 16-38.

Marekfia, W., Nissen, V., & für Dienstleistungen, F. W. (2012). Anforderungen an ein strategisches GRC-Management. In GI-Jahrestagung (pp. 731-745).

Okereke, C., & McDaniels, D. (2012). To what extent are EU steel companies susceptible to competitive loss due to climate policy. Energy Policy, 46, 203-215.

OMAR, O.; MENEGAZZO, C.; STEFANI, C.; MACEDO, M.; SELL, D. [anti]fragility of technological and innovation parks towards extreme events: an assessment and analysis model. In: VI Congreso Internacional de Conocimiento e Innovación, 2016, Bogotá. Anais do VI Congreso Internacional de Conocimiento e Innovación, 2016.

OMAR, O.; CUNHA, C. J. C. A.; SELL, D. Sistemas de Gestão de Conhecimento e Gestão commercial e de vendas: uma revisão systemática. In: VI Congreso Internacional de Conocimiento e Innovación, 2016, Bogotá. Anales del VI Congreso Internacional de Conocimiento e Innovación, 2016.

Pupke, D. (2008). Compliance and corporate performance: the impact of compliance coordination on corporate performance. BoD–Books on Demand.

Sampaio, R. F., & Mancini, M. C. (2007). Estudos de revisão sistemática: um guia para síntese criteriosa da evidência científica. Braz. J. Phys. Ther.(Impr.), 11(1), 83-89.

Saito, R., & Silveira, A. D. M. D. (2008). Governança corporativa: custos de agência e estrutura de propriedade. Revista de Administração de Empresas, 48(2), 79-86.

Silva, E. R. G. D. (2016). Arquitetura de conhecimento para e-participação. Tese de doutorado, EGC – UFSC.

Rother, E. T. (2007). Editorial: revisão sistemática x revisão narrativa. Acta Paulista de Enfermagem, 20(2), p. v-vi.

Taleb, N. N. (2012). Antifragile: Things that gain from disorder (Vol. 3). Random House Incorporated.

Wilhelm, M. M., Blome, C., Bhakoo, V., & Paulraj, A. (2016). Sustainability in multi-tier supply chains: Understanding the double agency role of the first-tier supplier. Journal of Operations Management, 41, 42-60.

# INVESTMENT DECISION IN THE AGENCY THEORY FRAMEWORK

**Ahmad Cahyo Nugroho,  Muhammad Firdaus,  Trias Andati dan Tony Irawan**
Ministry of Industry, Indonesia
Bogor Agricultural University, Indonesia
ac.nugroho999@gmail.com, mfirdaus@ipb.ac.id, tonyirawan82@gmail.com

**Abstract**. This studies aims to observe the development of literature on company investment decisions and to decide what research should be conducted further on company investment decisions in the theoretical framework of agency theory. The methods used were bibliometric network analysis and literature review. This study has mapped out the literature on company investment decisions based on agency theory. This study shows that the topics on competition research, corporate governance, and capital structure are closely related to the company investment decisions in the theory of agency, and it is worth investigating. Therefore, it is necessary to develop further empirical research related to company investment decisions in the framework related to agency theory by analyzing the influences of competition, corporate governance, and capital structure comprehensively on invesment decision.

**Keywords:** investment decision, bibliometric network analysis, capital structure, corporate governance, competition

**Abstrak.** Tujuan dari penelitian ini adalah untuk observasi pengembangan literatur pada keputusan investasi perusahaan dan dapat menentukan penelitian yang harus dilakukan selanjutnya mengenai keputusan investasi perusahaan dalam kerangka teori keagenan. Metode yang digunakan kali ini ialah *bibliometric network analysis* dan kajian literatur. Penelitian ini telah memetakan literatur mengenai keputusan investasi perusahaan. Kajian ini menunjukan bahwa topik penelitian persaingan, tata kelola perusahaan, dan struktur modal berkaitan erat terhadap keputusan investasi perusahaan dalam teori keagenan, dan layak untuk di teliti lebih lanjut. Oleh karena itu perlu dikembangkan penelitian empiris lebih lanjut terkait keputusan investasi perusahaan dalam kerangka terkait teori keagenan dengan menganalisis terkait pengaruh pengaruh persaingan, tata kelola perusahaan, dan struktur modal secara komprehensif.

**Kata kunci:** keputusan investasi, *bibliometric network analysis*, struktur modal, tata kelola perusahaan, persaingan

## INTRODUCTION

Decisions making regarding the utilizations of funds owned by a company, there are some differences in interests among the managers, employees, shareholder, and debtholder; similarly, the same thing happens when a corporate or a company will take investment decisions. On the other hand, investment decisions are important for the company in maintaining its sustainability and in growing its business. Investment decisions are important because they are related to the profitability that the company will earn, and this will even give impacts on the company value supported by Karuna (2007), Laksmana and Yang (2015), and Akdogu and MacKay (2012) who state that investment decisions can maximize the company value.

In addition, according to Gilbert and Lieberman (1987), investment would reduce the possibility of the competitor to expand, but temporarily it is influential. In

the development of studies related investment decisions has a close relationship with agency theory. This is because the agency theory is established based on the seven basic assumptions of personal interests, goal conflicts, restricted rationality, information asymmetry, efficiency excellence, risk aversion, and information as a commodity (Eisenhardt 1989). The studies which were concerned with the relationship between agency and the company investment decisions include Myers and Majluf (1984), Jensen (1986) and Fazzari *et al.,* (1988) who argue that the problem of information asymmetry between management and financial institutions and agency conflicts between controlling shareholders and minority investors and between management and shareholders have been proved to significantly affect company investment decisions. Meanwhile, according to Aivazian (2005), agency problems arise from the interaction among shareholders, debt-holders, and management.

According to Davis *et al.,* (1997), there are various mechanisms of agency theory in order to protect the interests of shareholders, minimize agency costs and ensure alignment interests. Jesen (1983) outlines two forms of agency theory that have been developed, namely, positivist and principal-agent. Positivist researchers have emphasized on governance mechanisms, especially in large corporations to identify agency issues and governance mechanisms that address agency problems. Eisenhardt (1989) states that positivist researchers have focused exclusively on special cases of principal-agent relationships between owners and managers of public companies. Meanwhile, the principal-agent research approach is a more general approach of agency theory where it explores the relationship between two parties such as workers and employers; legal consultants and their clients, and buyers and sellers. The principal-agent flow focuses on the technical and mathematical relationships of the specific details of contracts among the principal-agents. In other words, the focus of the principal-agent is to determine the optimal contract (Eisenhardt 1989). Positive accounting theory (Watts and Zimmerman 1986) proposes three hypotheses i.e. bonus plan hypothesis, debt/equity hypothesis, and political cost hypothesis which implicitly recognize three agency forms i.e. between owner and management, between creditor and management, and between government and management. Therefore, broadly speaking, the principal is not only the owner of the company, but it can also be a shareholder, creditor, or government.

The agency problem was initially explored by Ross (1973), while the detailed theoretical exploration of agency theory was first expressed by Jensen and Mecking (1976) stating that the manager of a company is the "agent" and the principal is the "shareholder". Jensen and Meckling (1976) describe the agency relationship in agency theory that a company is a nexus of contract between the owner of the economic resources (principal) and the manager who takes care of the utilization and controls the resources. According to Messier *et al.,* (2006), this agency relationship results in two problems: (a) the occurrence of information asymmetry, where management generally has more information on the actual financial position and the entity operation position of the owner; and (b) the occurrence of a conflict of interest due to inequality of goals, where management does not always act in the interests of the owner. In an effort to overcome or reduce the agency problem, agency cost which will be borne by both principal and agent must be provided. In regard to agency issues, corporate governance, which is a concept based on agency theory, is expected to serve as a tool to give assurance to investors that they will receive return on the funds they have invested.

Previous study showed that the company investment sensitivity on stock prices is determined by the extent to which the company has asymmetric information (Chen *et al.,* 2007) and problem agency (Jiang *et al.,* 2011). According to Jensen and Meckling (1976), asymmetric information makes investment decisions inefficient. The study by Guariglia and Yang (2015) documented strong evidence of investment inefficiency which was explained through a combination of financing constraints and agency problems. Two significant conclusions were emerged from this major finding. On the one hand, limited access to capital markets has caused many Chinese companies to be under-investment. In addition, weak corporate governance structures lead shareholders or managers to overinvest their free cash flow in projects with negative NPV. Shin and Kim (2002) also found that corporate investment decisions are influenced by the company agency cost, and this was supported by the study conducted by Hirt *et al.,* (2010). Based on the literature, it can be seen that the company investment decisions have relation closeness to the agency theory because in making investment decisions, they are inseparable from the conflicts among the stakeholders described in the agency theory.

**Problems.** There is limited information on the factors that determine investment decisions within the framework of agency theory. In addition, various studies show inconsistent results due to inefficient market conditions (Mutlu *et al.,* 2016; Young *et al.,* 2008; Bruce *et al.,* 2005; Sundaramurthy and Lewis 2003; Hill 1992), and agency theory applies only to the anglo saxon context (Bruce *et al.,* 2005). Hill (1992) who considers that in the event of a discrepancy in market conditions ignored by the researchers, the paradigm in agency theory should be comprehensively studied.

In the studies by Mutlu *et al.,* (2016) and Bruce *et al.,* (2005), there was some doubt to the application of agency theory on companies that develop in the context of non-anglo saxon culture and on companies that are still developing (Young *et al.,* 2008). In addition, Sundaramurthy and Lewis (2003) state that in an efficient market, agency theory will produce inconsistent results. This is in accordance with Hill (1992) that states that agency theory works based on the assumption that the market is efficient and adjusts for rapid changes. Inconsistent results in agency theory in previous studies led to the importance of prior mapping of the obtained literature, which can then be further studied.

**Objective.** The study aims to see the mapping of literature on research topics related to company investment decisions within the framework of agency theory and to find out research topics on company investment decisions within the framework of agency theory can be further developed.

## LITERATURE REVIEW

According Karuna (2007), product market competition of an industry influences managerial decisions. One of the managerial decisions is that investment decisions, based on Laksmana and Yang (2015) and Alimov (2014) show that competition also influences company investment decision. Competition itself is essential for growing a business because competition leads to business efficiency (Yi 2014), and Griffith (2001) also points out that product market competition plays an important role in reducing agency costs. In addition, competition allows companies to optimize the performance of their managers toward competitors (Laksmana and Yang 2015).

The study by Jiang *et al.,* (2015) found that high investment in high competition would increase the value proposition of the company. In addition, the study found a positive relationship between product market competition and corporate investment using a sample of Chinese manufacturing company in the period of 1999-2010. The study by Alimov (2014) also found that key corporate decisions are fundamentally influenced by product market competition. Theories and other empirical evidence also show that corporate investment decisions are influenced by competitive pressures, and competition as a different business cycle affects the amount and stability of the company cash flow (Mello and Wang 2012).

Corporate governance plays an important role in the management of a company because it will affect the performance of the company as stated by Byun *et al.,* (2012) that internal corporate governance has a positive effect on corporate value, and Ammann *et al.,* (2010) also shows that corporate governance will increase the value of the company. Similarly, investment decisions cannot be separated from the role of corporate governance (Guariglia and Yang 2015). Ammann *et al.,* (2013) indicates that corporate governance significantly increases the value of company in non-competitive industries only. In addition, the study found that good corporate governance for companies in non-competitive industries makes them have more capital expenditures, spend less on acquisitions, and tend to diversify. According to Zhou *et al.,* (2016), governance of the board of directors has a positive effect on managerial risk taking, indicating that board governance will result in higher investment in R & D expenditure and lower investment in capital expenditures, not only in the current year but also for the following years.

Giroud and Mueller (2011) found that weak corporate governance in competitive industries has lower returns on equity, poorer operating performance, and lower corporate value. The researchers also found that weak corporate governance in competitive industries is more likely to be targeted by hedge fund activists, suggesting that investors are taking action to reduce inefficiency.

The results of the study by Hu and Liu (2014) show that companies having CEOs with more diverse career experiences exhibit less cash flow sensitivity and exploit outside funds more, including bank loans and trade credit. Hossain *et al.,* (2000) found that the percentage of external directors is positively related to the investment opportunities of the company. Bathala and Rao (1995) and Hutchinson (2002) found a negative relationship between the proportion of external directors and the company growth rate. In contrast, Hossain *et al.,* (2000) found that the percentage of external directors is positively related to the investment opportunities of the company.

Anwar and Sun (2014), Aivazian *et al.,* (2003) examined the impact of financial leverage on corporate investment decisions using information in the Canadian public companies. The results also show that leverage is negatively related to investment, and this negative effect is significantly stronger for companies with low growth opportunities compared to those with high growth opportunities.

Anwar and Sun (2014) argue that an increase in foreign ownership, which refers to the level of foreign investment in domestic companies, can affect the leverage of the companies. The empirical results show that foreign investment has a negative influence and is significant statistically on the leverage of domestic companies in the manufacturing sector in China. By subsector, the results show that the impact of foreign investment on the leverage of domestic enterprises in China textile industry is negative

and significant, and this impact is much stronger than the overall impacts on the manufacturing sector. Meanwhile, the impact of foreign investment on leverage of domestic enterprises in the subsector of electrical machinery and equipment for the manufacturing industry is negative but less than the overall impacts in this type of industry.

## METHOD

This study used the literature review method on company investment decisions. The keyword data, titles and abstracts from 223 literatures obtained from the literature search engine were collected, and they were then analyzed using bibliometric network analysis. From this analysis, first we obtained some information on overlay visualization to see the development of literature obtained, and the next information was network visualization used to identify clusters and relationships among the research topics. Finally, density visualization was identified to see the density of the topics of the literature obtained. The results of bibliometric network analysis serve as a reference for a more in-depth literature review of the matters or factors affecting the company investment decisions. In the literature review, the selection and classification of the literature obtained based on the clusters from the results of the bibliometric network analysis was firstly conducted. Furthermore, a review on the research results of the study according to the clustering was carried out. Following this, the contradictions and differences of research results from the existing literatures were analyzed, and finally, the conclusions of the results of the analysis of the existing literature were made.

## RESULT AND DISCUSSION

### Bibliometric Network Analysis
This analysis viewed the development of the literature on investment decisions within the framework of agency theory using bibliometric network analysis of 223 literatures. The results obtained are as follows.

20

**Figure 1.** Overlay Visualization

Based on Figure 1, it can be seen that from the literature obtained, the research topics commonly discussed from 2002 to 2004 included the capital structure, leverage, debt, and firm value; on the other hand, from 2004 to 2006, the topics included information, risk, corporate investment, and R & D expenditure. Furthermore, from 2006 to 2008, the research topics covered manager, expenditure, agency cost, agency problem, and capital investment, and from 2008 to 2010, they covered shareholder, R & D investment, growth, product market competition, competition, profitability, ownership structure, agency theory, corporate finance, cash holding, firm performance, and corporate governance.



**Figure 2.** Density Visualization

21

Figure 2 shows the density of an item or topic of a study. The higher the number of topics around the point and the higher the weight of the topics surround them, the closer the point to the red color. Conversely, the smaller the number of items or topics, the lower the weight of the topics surround them, the closer the point to the blue color. Based on the density and proximity distance of the topics as shown in Figure 2, it can be seen that product market competition and competition have a large density and a close distance The topics that have a close distance and large density are corporate governance, shareholder, manager and agency cost supported by the topics of agency problem, cash holding, and corporate finance whose position is next to the location with less density. Furthermore, the larger density and close proximity to most research topics include expenditure with capital investment, information, and firm value around the area; moreover, risk, debt, and leverage with firm performance around them followed by growth with R & D investment, R & D expenditure, and corporate investment next to the area; lastly, capital investment with profitability, agency theory and ownership structure are located around the area.



**Figure 3.** Network Visualization

The figure shows that the research themes of investment, competition, corporate governance, capital structure, and agency theory have direct 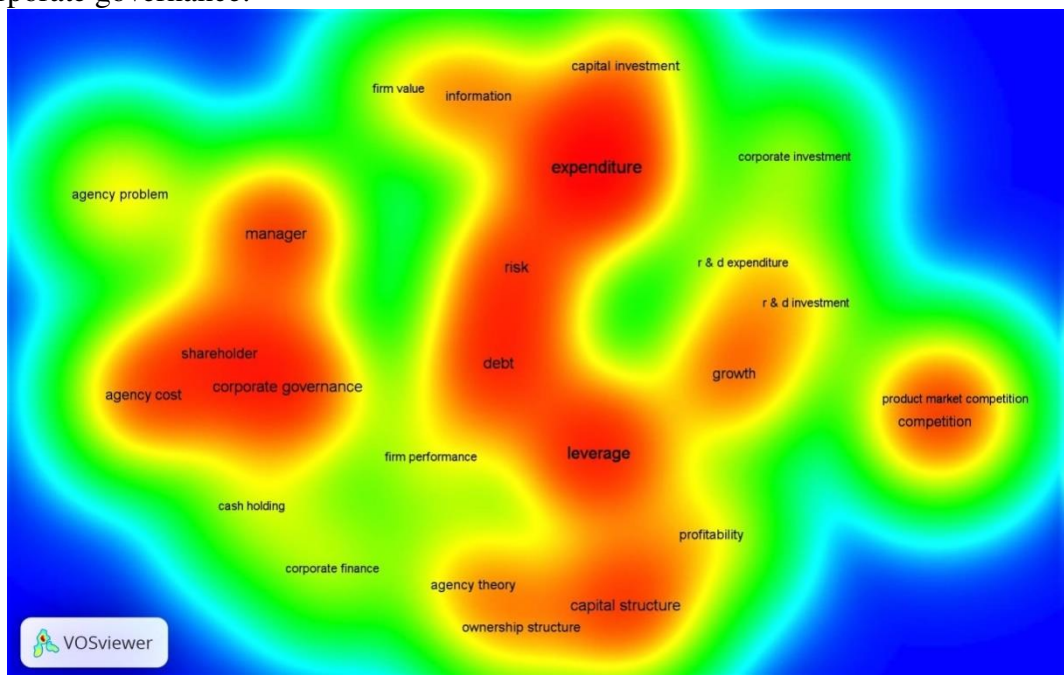and indirect relevance to each other. The theme of investment research is shown by the topics on capital investment, corporate investment, expenditure, R & D expenditure, R & D investment, and expenditure. Moreover, the theme of competition research is shown by the topics on product market and competition. Furthermore, the theme of capital structure research is shown with the topic of leverage, capital structure, and debt while the theme of corporate governance research is shown by research topics on corporate governance, ownership structure, manager, and share holder. Also, the theme of agency research is shown by agency theory, agency problem, agency cost, and information.

Based on Figure 3, there are three major clusters, and the first cluster marked by red color includes the research topics on capital investment, corporate investment, expenditure, R & D expenditure, R & D investment, growth, information, risk,

22

competition, and product market competition grouped into one cluster. The second cluster marked by a blue color includes the research topics on agency theory, agency problem, agency cost, corporate governance, manager, and cash holding. The third cluster in green color covers research topics on corporate finance, capital structure, debt, leverage, profitability, and ownership structure.

When viewed from Figure 3, research with the theme of competition (product market competition and competition) is correlated directly with the topics of capital structure, risk, and growth. As for the theme of capital structure (leverage, capital structure, and debt), it is directly related to the research topics on product market competition, profitability, growth, risk, ownership structure, corporate finance, firm performance, information, firm value, expenditure, capital investment, expenditure, R & D expenditure, manager, agency cost, agency problem, and agency theory. Furthermore, the theme of corporate governance research (corporate governance, ownership structure, manager, and shareholder) is directly related to the research topics on information, expenditure, R & D expenditure, firm performance, risk, leverage, debt, capital structure, cash holding, agency cost, agency Problem, and agency theory.

From the analysis of bibliometric network analysis, it can be concluded that the development of literature starts from the themes of capital structure, investment, agency theory, competition to corporate governance. The results of the analysis show that there are 3 clusters of capital structure, corporate governance, and investment and competition. The results of bibliometric network analysis also show the relevance of other research topics with investment decisions. Previous research reinforced the results of the bibliometric network analysis which show that investment decisions are influenced by factors such as competition (Jiang *et al.,* 2015; Griffith 2001; Alimov 2014; Flammer 2013; Laksmana and Yang 2015; Cheung 2012; Yi 2014; Fresard and Valta 2013; Mello and Wang 2012), corporate governance (Guariglia and Yang 2015; Francis *et al.,* 2013), and leverage (Anwar and Sun 2014; Aivazian *et al.,* 2003). (3 point of views or 3 cluster). From the bibliometric network analysis results supported by the existing literature, the next stage will be the literature study on company investment decisions from the 3 cluster i.e. cluster 1 product market competition, cluster 2 corporate governance, and cluster 3 capital structure.

**Cluster 1. Product Market Competition.** Laksmana and Yang (2015) contribute to the literature by providing evidence of the role of product market competition in disciplining management investment decisions. First, the results show that competition encourages managers to invest in risky investments. One potential explanation for this result is that competition reduces the opportunity for resource shifting for management personal benefits. Another explanation is that the power of competition affects management to take on more risks for long-term company survival. Second, these results indicate that competition makes management more disciplined on cash usage. Overall, research results provide support for corporate governance functions from product market competition to corporate investments and show that companies in more competitive industries take more risks as measured by capital expenditure, R & D expenditure, and standard deviations from stock returns than those in less competitive industries.

According to Chen *et al.,* (2014), product market competition can lower equity capital cost, indicating that the agency cost competition will become more efficient. Lin *et al.,* (2012) state that increased competition of product market will improve efficiency

of banking and financing in labor intensive industries. The research by Jiang *et al.,* (2015) also shows that there is a positive relationship between competition of product market and corporate investment in a country with a strong and predictable economy like China.

On the contrary, Grullon and Michaely (2008) had somewhat a different study from other studies i.e. examining the interaction between product market competition and manager decisions to distribute cash to shareholders. The study used the Herfindahl-Hirschman Index (HHI) from the census of producers as a proxy for product market competition and used a sample of large manufacturing companies, and it is found that companies in less competitive industries have a significantly lower payout ratio than those in more competitive markets. Overall, these findings are consistent with the idea that the power of competition encourages managers to be disciplined in paying for excess cash and with the idea that dividends are the "outcome" of the external factors. Moreover, Baggs and Bettignies (2007) state that competition has a direct pressure, significant effect, and a significant institutional effect on managerial efficiency. These effects make the company understand the importance of quality improvement.

Cheung (2012) investigated how competition in the product market together affected corporate investment and financing behavior during the period of 1996-2008. The research explained the role of competition in different market structures. The results found that higher market product competition generally resulted in higher reduction in investment. However, these results do not apply to concentrated industries where competition of product market leads to significant investment expenditure. Meanwhile, Yi's (2014) study discusses the impact of product market competition on the investment efficiency of the company from both sides. First, in terms of corporate governance, product market competition facilitates institutional investors in disciplining companies, and this makes investment more efficient. However, if it focuses on the production information, product market competition reduces the company incentive to acquire information. Companies in a competitive industry receive different signals on exact demand which leads to less efficient investment.

Grosfeld and Tressel (2001) study show that product market competition has a positive and significant impact on performance. Finally, product market competition and good corporate governance tend to reinforce one another instead of being a substitute. Competition has no significant impact on performance for companies with bad corporate governance; in contrast, it has significant positive impact on companies with good corporate governance.

In general, with the competition, agency problems and agency costs will decrease, and this benefits the company because the competition will improve the performance of management and corporate governance. However, there is still inconsistency from research on product market competition such as the study by Cheung (2012) which shows inconsistency in the influence of competition on investment, i.e. if market structure is different, it will have a difference influence on investment. Furthermore, Laksmana and Yang (2015) argue that competition makes management more disciplined on the use of cash, but they also state that competition encourages managers to make risky investment decisions. As it is in the study of Yi (2004) which shows that competitive industries are less efficient in corporate

investment. In addition, lack of competitive research on investment in the framework of agency theory in Indonesia makes this necessary to be further studied.

**Cluster 2. Corporate Governance.** The study of Chen *et al.,* (2014) that examined the level and type of ownership with the allocation of investment funds found that ownership of a company has an important role in determining investment behavior of the company. Based on the study by Francis *et al.,* (2013), corporate governance is a key determinant of investment sensitivity to internal cash flows; in addition, the level of corporate governance has a positive effect on access to funding. The results of this study support the argument of agency theory which suggests the existence of incentives and monitor costs of managerial actions can lead to friction funding and affect the behavior of corporate investments. As it is in the research of Shin and Kim (2002), company investment decisions are also influenced by agency costs.

According to Chen *et al.,* (2015), companies that have more concentrated shares/stocks have a drive to conduct over-investment while those with higher stock trading proportion, board size of supervisors, and leverage will mitigate over-investment. For underinvestment firms, their evidence shows that firms with higher state-ownership concentration, larger board size of directors or higher proportion of outside directors are associated with severer under-investment, while firms with higher leverage or higher proportion of tradable shares alleviate under-investment.

According to Chen *et al.,* (2013), foreign institutional shareholdings make a difference in the level of asymmetry information, and problem agency makes a difference in investment behavior. The shareholdings will mitigate the problem agency and asymmetric information by improving corporate governance and transparent financial management. The study by Almeida and Dalmacio (2015) investigated how interaction between product market competition and corporate governance improves analyst accuracy estimates and reduces estimates of deviations. The study used a sample of public companies in Brazil. The results show that competitive industries provide incentives to improve the flow of information but not necessarily to improve quality. However, good corporate governance improves the financial reporting process, and consequently, the quality of analyst estimates becomes more accurate. The main evidence of the study is that analysts covering companies in competitive industries with strong corporate governance are the most accurate ones. Similarly, according to Guadalupe and Gonzalez (2010), competition policy can have an important influence on corporate governance.

The contract structure is influenced by the corporate governance. With the involvement of foreign investors in ownership, a better management system will be created due to the current knowledge and technologies that enter the company (Guariglia and Yang 2015; Chen *et al.,* 2014; Shin and Kim 2002). In addition, managerial share ownership mechanisms can improve the performance of the company and shareholder values (Jensen and Meckling 1976, Morck *et al.,* 1988; Wright *et al.,* 1996; Lins 2003; Wei *et al.,* 2005). In the agency theory, the option of the offering of the managerial shareholding interestingly encourages managers to improve their company performances (Sundaramurthy and Lewis 2003).

According to Connelly *et al.,* (2017), dividend payout initiation is associated with stronger governance (stronger positions of the shareholders and independent board). The managerial shareholding by the CEOs has a positive relationship with the dividend initiation at companies whose governance is highly strong. The study found

25

that when initiation was due to significantly stronger governance, the governance was then related to corporate investment opportunities, while for weaker corporate governance, the relationship was not observed.

Based on Chen *et al.,* (2017), the shareholdings of the government and foreign institutions are associated with different levels of asymmetric information and problem agencies. The studies found strong evidence that government shareholding have undermined the sensitivity of investment opportunities, thus increasing investment inefficiency. Conversely, the foreign shareholding has strengthened the sensitivity of investment opportunities, thereby increasing investment efficiency. In addition, I found that the relationship between foreign ownership and investment efficiency is stronger as the government relinquished control and the governance level of the government institutions is weaker.

According to Mykhayliv and Zauner (2017), the majority of shareholders and the increase in state share ownership have a negative influence on the company investment, such as domestic share ownership and financial institutions whereas the insider shareholdings and industry and financial groups do not affect investment.

Based on the existing literature, it can be concluded that in general, corporate governance is influential on company investment decisions; however, there are differences in the influences due to the level of competition or the structure of market competition that makes the governance different. In addition, the existing corporate governance studies still have shortcomings as in the study of Connelly *et al.,* (2017), who studied the relationship between weak corporate governance and unobserved investment decisions. Moreover, there is still inconsistency in terms of influence of share ownership on investment decisions as Mykhayliv and Zauner (2017) state that insider shareholdings and industrial and financial groups do not influence investment. It is, therefore, necessary to further study the influence of corporate governance on more specific investment decisions.

**Cluster 3. Capital Structure.** According to Firth *et al.,* (2012) in addition to shareholder, there are also debt holders who will monitor the company performance and business decisions of the company including investment decisions due to the leverage. The existence of problem agency in the research by Aivazian (2005) is "overinvestment" due to conflict between management and shareholders. Managers have a tendency to expand in the company scale only and pay less attention to future corporate value after investments; moreover, they even take on poor investment projects and reduce shareholder wealth. Management ability to implement the policy is limited by the availability of cash flow, and this constraint can be further reduced through debt financing. By issuing debt, the company has to pay interest and principals that put pressure on management not to allocate funds for poor investment projects. Similarly, Anwar and Sun (2014), Aivazian et al. (2003), Guney *et al.,* (2011) in their studies state that leverage affects investment decisions.

In order to prevent over-investment and increase control over corporate management, capital structure policy is often used. With the interest cost and creditor supervision, the company management will be more careful in investing for the company (Anwar and Sun 2014, Guney *et al.,* 2011 Aivazian *et al.,* 2003). Guariglia and Yang (2015) argue that the limited access to capital markets causes many Chinese companies to be under-investment. In addition, weak corporate governance structure

leads shareholder managers or controllers to overinvest their free cash flow in projects with negative NPV.

The findings obtained by Moenadin *et al.,* (2013) indicate that there is a significant relationship between selected industrial capital structure and product market competition. Similarly, the study by Pandey (2004) shows that there is a relationship between capital structure and market forces due to the complex interactions of market conditions, agency problems and bankruptcy costs. According to Rathinasamy (2000), companies that have more monopoly power use longer-term financing and debt. Based on the research by Guney *et al.,* (2011), Chinese companies tend to adjust their leverage ratios. There is a relationship between the intensity of competition and the leverage ratio, which supports the theory of predation.

According to Munisi (2017) who studied the determinants of capital structure in Sub-Saharan Africa, capital structure is negatively related to profitability and tangible assets, supported by pecking-order theory and trade-off theory. The findings show that capital structure is positively associated with free cash and corporate growth, consistent with agency theory and pecking-order theory.

Nevertheless, there are also different study results on agency theory in capital structure as in Banga *et al.,* (2017) who state that the determination of capital structure in the small-medium scaled businesses in India is influenced by the application of pecking order theory and trade off theory, but there is no evidence for the application of agency theory. Also, Huang *et al.,* (2016) examined the capital structure of the small-medium scaled businesses in China, and the results show that leverage is influenced by executive shareholders and excess cash compensation. However, institutional share ownership does not affect leverage level that is more influenced by traditional factors such as taxes and operating cash flows.

Empirical results from Dawar's (2014) research in India show that leverage has a negative effect on the financial performance of Indian companies, which contradicts the assumption of agency theory as it is generally accepted in developed and developing countries. Consequently, agency theory postulate should be viewed with a different perspective in India by considering the nature of bond market and dominance of state-owned banks in providing loans to the corporate sector.

Based on the existing literature, it can be concluded that with leverage, there is a conflict among shareholders, management and debtholders that can be seen from the agency theory. The sub-sectoral differences, foreign share ownership, and market structure will affect the company capital structure, and Pecking order theory and trade off theory affect the capital structure. On the contrary, the agency theory on capital structure is still inconsistent because a number of studies show that the agency theory is not found in determining capital structure.

**The implication of research development.** According to Hill (1992), there are paradigms in analyzing agency theory, namely, corporate strategic behavior; principal-agent contract structure; monitoring and strengthening of principal-agent contracts, and evolutionary processes that alter the structure of the principal-agent contract and the institutional structure that issues the contract. From the study, it was found that differences from market conditions resulted in inconsistent results, and by analyzing them from various paradigms, more comprehensive results in the application of agency theory can be obtained. Based on the literature review on investment decisions within

27

the framework of agency theory associated with the topics of product market competition research, corporate governance, and capital structure, inconsistent results still exist.

Therefore, further research is required to provide a more comprehensive review related to the company strategic behavior paradigm by using product market competition indicators affecting the company environmental conditions (Jiang *et al.,* 2015; Griffith 2001; Alimov 2014; Flammer 2013; Laksmana and Yang 2015; Cheung 2012; Yi 2014; Fresard and Valta 2013; Melo and Wang 2012); The principal-agent contract structure is indicated by implementation of good corporate governance within the company (Guariglia and Yang 2015; Chen *et al.,* 2014; Shin and Kim 2002); and the monitoring and strengthening of the principal-agent contract use the capital structure policy indicator, and with the existing of the interest cost and debtholders supervision, the company management will be more careful in investing the company (Anwar and Sun 2014, Guney *et al.,* 2011 Aivazian *et al.,* 2003). From this studies, it is expected to provide a study that can provide consistent results in inefficient market conditions.

## CONCLUSION

This study has mapped out the existing literature using bibliometric network analysis that can show that the development of literature starts from the themes of capital structure, investment, agency theory, competition and corporate governance. Then the results of the analysis indicate that there is a linkage of topics directly or indirectly, and this analysis shows from the existing literature that there are 3 major clusters of capital structure, corporate governance, and competition.

Based on the existing literature review, the research topics of product market competition, corporate governance, and capital structure are research topics that have closeness relation to company investment decisions within the agency theory framework. The previous studies generally stated that product market competition, corporate governance and capital structure affect investment decisions taken by the company and can be discussed with the agency theory. However, there are still deficiencies and inconsistencies in the results of the previous studies, and there is still lack of research on the influence of product market competition, corporate governance, and capital structure on corporate investment decisions within the framework of the agency theory that was thoroughly studied. Therefore, further empirical studies on the impact of product market competition, corporate governance, and capital structure on investment decisions in comprehensive agency theory are required.

## REFERENCE

Aivazian V, Ge Y, Qiu J. (2003). "The impact of leverage on firm investment: Canadian evidence". *Journal of Corporate Finance*. 11(1-2): 277-291. Doi:10.1016/S0929-1199(03)00062-2

Akdogu E, Mackay P. (2012). "Product markets and corporate investment: theory and evidence". *Journal of Banking & Finance*. 36 (2): 439–453.

Alimov A. (2014). "Product market competition and the value of corporate cash: evidence from trade liberalization". *Journal of Corporate Finance*. 25: 122–139.

Almeida JEF, Dalmácio FZ. (2015). "The effects of corporate governance and product market competition on analysts'forecasts: evidence from the Brazilian capital market". *The International Journal of Accounting*. 50(3): 316-339.

Ammann M, Oesch D, Schmid MM. (2013). "Product market competition, corporate governance, and firm value: evidence from the EU-Area". *European Financial Management*. 19(3): 452-469. doi: 10.1111/j.1468-036X.2010.00605.x

Ammann M, Oesch D, Schmid MM. (2013). "Product market competition, corporate governance, and firm value: Evidence from the EU area". *European Financial Management*. 19(3): 452-469.

Anwar S, Sun S. (2014). "Can the presence of foreign investment affect the capital structure of domestic firms?". *Journal of Corporate Finance*. 30: 32-43. doi: 10.1016/j.jcorpfin.2014.11.003

Baggs J, Bettignies JE. (2007). "Product market competition and agency costs". *The Journal of Industrial Economics*.  55(2): 289-323.

Banga C, Gupta A. (2017). "Effect of firm characteristics on capital structure decisions of Indian SMEs". *International Journal of Applied Business and Economic Research*. 15(10): 281-301.

Bathala CT, Rao RP. (1995). "The determinants of board composition: an agency theory perspective". *Managerial and Decision Economics*. 16: 59–69.

Bruce A, Buck T, Main B. (2005). "Top executive remuneration: A view from Europe". *Journal of Management Studies*. 42(7): 1493-1506.

Byun HS, Lee JH, Park KS. (2012). "How does product market competition interact with internal corporate governance?: evidence from the Korean economy". *Asia-Pacific Journal of Financial Studies*. 41(4): 377-423 doi:10.1111/j.2041-6156.2012.01077

Chen C, Li L, Ma MLZ. (2014). "Product market competition and the cost of equity capital: evidence from China". *Asia-Pacific Journal of Accounting & Economics*. 21(3): 227–261.

Chen Q, Goldstein I, Jiang W. (2007). "Price informativeness and investment sensitivity to stock price". *Review of Financial Studies.* 20: 619-650.

Chen R, Ghoul AE, Omrani G, Wang H. (2013). "Do state and foreign ownership affect investment efficiency? Evidence from Privatizations". 42: 408-421. doi: 10.1016/j.jcorpfin.2014.09.001

Chen X, Sun Y, Xu X. (2015). "Free Cash Flow, Over-Investment and Corporate Governance in China". *Pacific-Basin Finance Journal*. 37: 81-103. doi: 10.1016/j.pacfin.2015.06.003

Chen R, El Ghoul S, Guedhami O, Wang H. (2017). "Do state and foreign ownership affect investment efficiency? Evidence from privatizations". *Journal of Corporate Finance*. 42: 408-421.

Cheung A. (2012). "Product market competition, Corporate Investment and financing: evidence from cash flow shortfall". Curtin University, School of Economics and Finance. Bentley Campus, Perth, Australia.

Connelly BL, Shi W, Zyung J. (2017). "Managerial response to constitutional constraints on shareholder power. Strategic Management Journal". 38(7): 1499-1517.

Davis JH, Schoorman FD, Donaldson L. (1997). "Toward a stewardship theory of management". *Academy of Management Review*. 22(1): 20-47.

29

Dawar V. (2014). "Agency theory, capital structure and firm performance: some Indian evidence". *Managerial Finance*. 40(12): 1190-1206.

Eisenhardt KM. (1989). "Agency theory: an assessment and review". *Academy of Management Review*. 14 (1): 57-74.

Fazzari SM, Hubbard RG, Petersen B. (1988). "Financing constraints and corporate investment". *Brookings Papers on Economic Activity*. 19: 141– 195.

Firth M, Malatesta PH, Xin Q, Xu L. (2012). "Corporate investment, government control, and financing channels: evidence from China's listed companies". *Journal of Corporat Finance*. 18: 433–450.

Flammer C. (2013). "Does product market competition foster corporate social responsibility?". *Strategic Management Journal*. 36(10): 1469-1485 doi: 10.1002/smj.2307

Francis B, Hasan I, Song L, Waisma M. (2013). "Corporate governance and investment-cashflow sensitivity: Evidence from emerging markets". *Emerging Market Review*. 15: 57-71.

Fresard L, Valta P. (2013). "Competitive pressure and Corporate Investment: Evidence from Trade Liberalization". University of Maryland

Gilbert RJ, Lieberman M. (1987). "Investment and coordination in oligopolistic industries". *Journal of Economics*. 18(1): 17-33.

Giroud X, Mueller HM, (2011). "Corporate governance, product market competition, and equity prices". *Journal of Finance*. 66: 563–600.

Griffith R. (2001). *Product market competition, efficiency and agency costs: an empirical analysis*. The institute for fiscal studies.

Grosfeld I, Tressel T. (2001). "Competition, corporate governance: substitutes or complements? evidence from the Warsaw stock exchange". *Economics of Transition*. 10(13): 525-551. doi: 10.1111/1468-0351.t01-1-00124

Grullon G, Michaely R. (2008). "Corporate payout policy and product market competition". *New Orleans Meetings Paper*. http://dx.doi.org/10.2139/ssrn.972221

Guadalupe M, Gonzalez FP. (2010). "Competition and private benefits of control". *Chicago Meetings Paper*. http://dx.doi.org/10.2139/ssrn.890814

Guariglia A. (2008). "Internal financial constraints, external financial constraints, and investment choice: evidence from a panel of UK firms". *Journal of Banking Finance* 32: 1795–1809.

Guney Y, Li L, Fairchild R. (2011). "The relationship between product market competition and capital structure in Chinese listed firm". *International Review of Financial Analysis*. 20: 41-51

Hill CWL, Jones TM. (1992). "Stakeholder – agency theory". *Journal of Management Studies.* 29(2): 131 – 154.

Hirth S, Homburg MU. (2010). "Investment timing, liquidity, and agency costs of debt". *Journal of Corporate Finance*. 16(2): 243–258

Hossain M, Cahan SF, Adams MB, (2000). "The investment opportunity set and the voluntary use of outside directors: New Zealand evidence". Accounting and Business Research 30(4): 263– 273.

Hu C, Liu YJ. (2014). "Valuing Diversity: CEOs' Career Experiences and Corporate Investment". *Journal of Corporate Finance*. 30: 11-31. doi: 10.1016/j.jcorpfin.2014.08.001

30

Huang W, Boateng A, Newman A. (2016). "Capital structure of Chinese listed SMEs: an agency theory perspective". *Small Business Economics*. 47(2): 535-550.

Hutchinson M. (2002). "An analysis of the association between firms' investment opportunities, board composition, and firm performance". *Asia-Pacific Journal of Accounting and Economics*. 9: 17– 39

Jensen MC, Meckling WH. (1976). "Theory of the firm: managerial behavior, agency cost and ownership structure". Journal of Financial Economics. 3: 305– 360.

Jensen MC. (1983). "Organization theory and methodology". *Accounting Review*. 56(2): 319-338.

Jiang F, Kim KA, Nofsinger JR, Zhu B. (2015). "Product market competition and corporate investment: evidence from China". *Journal of Corporate Finance*. 35: 19-210. *doi: 10.1016/J.Jcorpfin.2015.09.004.*

Jiang F, Kim K.A. (2015). "Corporate governance in China: a modern perspective". *J. Corp. Finance* 32: 190-216.

Jiang L, Kim JB, Pang L. (2011). "Control-ownership wedge and investment sensitivity to stock price". *Journal of Banking & Finance* 35: 2856-2867.

Karuna C. (2007). "Industry product market competition and managerial Incentives". *Journal of Accounting and Economics*. 43 (2-3): 275–297.

Laksmana I, Yang Y. (2015). "Product market competition and corporate investment decisions". *Review of Accounting and Finance*. 14(2): 128-148. http://dx.doi.org/10.1108/RAF-11-2013-0123.

Lin JY, Sun X, Wu HX. (2012). "Banking structure, labor intensity, and industrial growth: Evidence from China". Working Paper: World Bank.

Lins KV. (2003). "Equity ownership and firm value in emerging markets". *Journal of Financial and Quantitative Analysis,* 38: 159-183.

Mello AS, Wang M. (2012). "Globalization, product market competition, and corporate investment". Working Paper. University of Wisconsin.

Messier WF, Glover SM, Prawitt DF. (2006). *Auditing and Assurance Services A Systematic Approach*. Fourth Edition. McGrw-Hill Companies,Inc. New York.

Moenadin M, Nayebzadeh S, Ghasemi M. (2013). "The relationship between product market competition and capital structure of the selected industries of the Tehran stock exchange". *International Journal of Academic Research in Accounting, Finance and Management Sciences*. 3(3): 221–2.

Morck R, Shleifer A, Vishny RW. (1988). "Management ownership and market valuation: An empirical analysis". *Journal of Financial Economics*. 20: 293-315.

Munisi GH. (2017). "Determinants of capital structure: Evidence from Sub-Saharan Africa". *International Journal of Managerial and Financial Accounting*. 9(2): 182-199.

Mutlu CC, Peng MW, van Essen M, Saleh S. (2016). "Agency theory and corporate governance in China: a meta-analysis". *Working Paper.* [Internet]. [diunduh 2016 Desember 10]. Tersedia pada: https://works.bepress.com/canan-mutlu/5/

Myers SC, Majluf NS. (1984). "Corporate financing and investment decisions when firms have information that investors do not have". *Journal of Financial Economics*. 13: 187–221.

Mykhayliv D, Zauner, KG. (2017). "The impact of equity ownership groups on investment: Evidence from Ukraine". *Economic Modelling*. 64: 20-25.

31

Pandey IM. (2004). "Capital Structure, Profitability and Market Structure: Evidence from Malaysia". *Asia Pacific Journal of Economics & Business*. 8(2): 78-91.

Rathinasamy RS, Krishnaswamy CR, Mantripragada KG. (2000). "Capital structure and product market interaction: an international perspective". *Global Business and Finance Review*. 5(2): 51−63.

Ross SA. (1973). "The economic theory of agency: the principal's problem". *American Economic Review*. 63(2): 134–139.

Shin HH, Kim YH. (2002). "Agency costs and efficiency of business capital investment: evidence from quarterly capital expenditures". *Journal of Corporate Finance* . 8(2): 139–158.

Sundaramurthy C, Lewis M. (2003). "Control and collaboration: Paradoxes of governance". *Academy of Management Review*. 28(3): 397-415.

Watts RL. Zimmerman JL. (1986). "Positive accounting theory : A ten a year perspective". *The accounting review*. 65(61): 131-156.

Wei KCJ, Zhang Y. (2008). "Ownership structure, cash flow, and capital investment: Evidence from East Asian economies before the financial crisis". *Journal of Corporate Finance.* 14 (2): 118–132

Wright P, Ferris SP, Sarin A, Awasthi V. (1996). "Impact of corporate insider, blockholder, and institutional equity ownership on firm risk taking". *Academy of Management Journal*, 39(2): 441-463.

Yi L. (2014). "Product market competition and investment Efficiency". *Tesis*. Hongkong (SE) : The University Of Hong Kong

Young M, Peng MW, Ahlstrom D, Bruton GD, Jiang Y. (2008). "Corporate governance in emerging economies: A review of the principal-principal perspective". *Journal of Management Studies,* 45(1): 196-220.

Zhou T, Li WA. (2016). "Board governance and managerial risk taking: Dynamic analysis". *Chinese Economy*. 49(2): 60-80.

32

# Bibliometric Analysis: Agency Theory in Accounting

**Isaac Francis Antwi**[a]

[a] University of Aveiro, Department of Economics, Management, Industrial Engineering and Tourism, 351-920276713 Email: isaacantwi@yahoo.com OR isaacantwi@ua.pt

**Abstract**

**Purpose:** This paper conducts a general bibliometric analysis review of agency theory in accounting (financial and management). The bibliometric analysis offers historical information on-trend and performance research.

**Methodology:** The study investigated the related literature in the agency theory and accounting (financial and management) from 1999-2019, obtained from the Scopus database. The literature-based documents are on the study of the scientific output and distribution of subject categories and journals. Keywords of the authors also have focused on determining the study hotspots.

**Findings:** The findings of this study show that annual production has increased over the period under investigation. The Critical Perspective on Account is the leading prolific journal and Accounting, Auditing and Accountability is a most influential journal. The result also shows that many top institutions are from the United Kingdom. Simultaneously, the United States of America leads the highest production and cited documents of related scientific articles.

**Originality /Value:** This study contributes on the awareness of using bibliometric analysis study to explore development in the scientific field, that is, the use of keywords to extract information for research growth in terms of the number of production and citations.

## 1.0 Introduction

Agency theory has gained its prominence in international organisations, academics, professional practices, and corporate bodies over the years now. The approach is base upon the principal-agent framework. Jensen and Meckling (1976) described agency relationship as a contract under which one or more persons (principal) engage another person (agent) to perform some services on their behalf, which involves delegating some decision-making authority to the agent. An agency relationship arises when a provider of funds appoints another to manage his interest. Proponents of agency theory believe that there is a tendency for agents, when left unmonitored, to engage in self-interest activities to the principal's detriment. The agency theory has also related to how the information is for external users (stakeholders), especially shareholders in an organisation, and internal users of information (management) of financial data like financial reports on a firm's performance, budget, and performance evaluation.

Over the years, several writers have provided a broader array of accounting overviews. Some used bibliometric measures to determine the general condition of the research field (Brown, 1996). Several others have also researched various essential aspects, like journal rankings (Qu et al., 2009). Moreover, others have concentrated on comparing accounting with other related disciplines. Nevertheless, none of them gave a full view of the current state of the art, taking into account all the modern instruments available to reflect a field with bibliometric indicators (Podsakoff et al., 2009). Considering the relevance and expansion of agency theory research and following the trend of researchers' concentration, which of interest to the masses on the study of the scientific production.

This paper aims to present a general bibliometric analysis of agency theory integrating into accounting (financial and managerial) research over a period from 1999-2019 with a new approach that combines several tools for representing the importance of bibliographic material found in Scopus database. This database is usually regarded as the most influential in academic research because it includes journals, articles, authors, years, and recognised the highest citation counts. The

number of publications, citation count analysis, and the *h*-index is currently a measure for representing the quality of a set of papers (Hirsch, 2005). It assumed that the number of documents indicates total production (TP), and the number of total citation (TC) count is the most influential research area. The h-index consider the two. Most of the results follow the general understanding of the direction of the scientific field analysed from Scopus. This article novelty is to use extant accounting literature and bibliometric analysis indicators to present the new trend of academic research.

The study analyses agency theory in accounting from Scopus. Some review papers have also focused on financial accounting and management accounting. Still, there is little research that analyses agency theory in financial and management accounting using bibliometric analysis. Bibliometric analysis is a useful tool evaluating countries, institutions, authors, journal authors (Tang et al., 2018). This study aims to achieve the following objectives by using bibliometric indicators:

1. Identify the top 10 trend players in terms of leading journals(production and citation), years(production), highly impactful articles, influential authors (production and citation), authors in a leading journal(production and citation), impactful institutions, country (production and citation) on the scientific field in the period under review.

2. Determine keywords that have drawn the attention of scholars subject area in the past, present, future direction.

According to Curty and Boccato (2005), keywords represent the words of the paper's subject under review. The choosing of the keywords used in this research is agency theory, financial, management, and accounting from Scopus.

The significant contributions of the research are the highlights as follows:

1. The bibliometric study explored recent 10 top development in the scientific field, using the bibliometric indicators. It will serve as the starting point for academia to research in unexplored areas yet potentially significant.

2. From this study, researchers can understand the inner structure of conducting research using bibliometric indicators to get a broad picture of this area.

3. The study will cover bibliometric analysis research growth regarding number production and total citation received over the years.

The study uses a bibliography analysis to answer the research questions set out above: reviewing existing literature to recognise main trends and problems, and proposing the justification for single-source reference for scholars and organisation management interested in bibliometric methods. The idea suggests workflow guidelines to carry out bibliometric studies in the future. This study provides a statistical survey of published papers and citations for calculating the effect on the field.

Note that the review is not limited to any language; however, English speaking countries dominated most on the publications.

The organisation of this work is as follows: We provide a literature review of agency theory and integrating financial accounting, and management; the methodology of the study; bibliometric analysis results; and conclusion.

## 2. 0 Literature Overview

### *2.1 Integrating agency theory in accounting*

The Agency theory arose from various authors', but the first author Adams Smith was the belief brain behind it in 1776 after which other authors took the inspiration from till now. Smith views that an organisation managed by a group of people or an individual who is not the real owners of the business, their chance of the business not being managed well will surface at the owner's expense. The argument is that the manager being an agent might use the company's property for his gain, which may bring a problem. Berle and Means (1932) found the research on agency theory, that, agent appointed by owners, and the agent might use the firm's property for his ends, which will create conflict between the principal and agent. Jensen and Meckling (1976), an agency relationship is also a contract between the principal and agent. Both parties work for their self-interest that leads to the agency conflict. They consider the agreement as a legal binding document between the principal and the agent principle uses a monitoring tool to curb the agent's activities with the view to

control cost. Ross (1973) identify the agency problem resulting from the decision on the compensation from the contract emanated from the firm and society.

The agency theory is to justify accounting research since the idea of approach provides a framework for explaining contractual relationships between a principal and an agent (Chi, 1989). In agency theory, an individual is motivated by self-interest, which could lead to the agency problem. It is useful to provide an overview of how agency theory should integrate into the financial and management accounting, and understanding how agency relationships work and the accounting information system.

An Agency theory argued that in a modern company, the distinction between the owners and management has led to disputes, where the agent (management) acquires more knowledge and thus appears to behave to their advantage, rather than to meet the desires of the principal (shareholders) (Berle and Means, 1932). Accounting information provides information to interested parties about the financial situation and results through a sound accounting system, which produces financial accounting reports and firms performance. Therefore, there is the need to share information about the company's financial position and performance that might be relevant and useful for their decision making through sound accounting information system (Brown et al., 2011). Schoenfeld (2020) examined large-scale data set of contracts using block investments from 1996 to 2018 on how the theory predict bilateral agreements. Specific contracts which include financial reporting and access to information resulted that contract arrangements relating to accounting information are significantly associated with information asymmetry measures between managers and shareholders. Cunha et al. (2016) investigate the accounting information's effect on the Portuguese mayor's re-election. The study considers agency theory as a starting point; the result shows that specific accounting information influences the re-election of mayors in Portugal, namely the financial accounting components. Similarly, Hiebl (2015) explores the various pay attitudes of Chief Financial Officers (CFOs) an (agents) which may align with the theory of agency. The result shows that CFOs who reports financial accounting to owners

expected more power in the owners' hands. This indicates that company owners' actions will affect and change the attitude of a company managers.

Management accounting is part of the accounting system that concerns measurement and information within an organisation, seeking to evaluate past decisions and improve future decision. Agency theory has been one of the approaches which integration has been inclusive in management accounting to incorporate conflict of interest, incentives problems and mechanism for adopting incentive package to control the agency problems. Such managerial accounting on individual interaction within a firm is essential to mitigate the underline agency problem by absorbing it into the employment contract. Brink et al. (2017) examine the effect of financial incentive to managers believing that such incentives will cause managers to engage on excessive risk-taking that effect the managers pay and another participant, the agent. The result indicates that the managers' incentives package is aligned with the organisation's interest to reduce the agency problem. Baiman et al. (2010) examine the effect of management accounting potential effects on firm performance on the informativeness and incentive performance evaluation. Finding suggests agent is paid based on the output of his contract workstation.

## 3.0 Methods

The study used bibliometric analysis as a combination of statistical methods and literature (Pritchard, 1969). Bibliometric analysis entails various steps; performance analysis (including data extraction, processing, analysis) and science mapping (networking and visualisation) (Cobo et al., 2011). The performance analysis seeks to assess individual and institution research and publication, whereas science mapping intends to show the structure and dynamic forces of science field. Bibliometric methods introduce quantitative evenness into the subjective evaluation of literature. The process provides evidence of theoretically developed classifications in a review paper. The implementation determines the effect frequency of the influence of research; the relationship indicators measure the connection between researchers and their different fields. The result is an overview of the research effort development, evolution, and quality (Ramos-Rodrıguez and Ruız-Navarro, 2004).

## *3.1 Data Extraction*

Data were extracted from the Scopus database as is considered to be one of the largest databases, extracting 348 articles in the related scientific field or business-related field. The team of words' agency theory', 'finance', 'management', and 'accounting' formed to recover the related documents, even though scholars may have used different texts to retrieve the records from search repository from 2010 to 2019 as at 27/02/2020.

## *3.2 Data Analysis*

The current study employs a combination of process among a lot, such as the total number of publications, the total number of citation counts and the h-index. Merigo et al. (2015) believe three most practical papers that define a group's value are the number of publications, citations count and h-index. The publication count has received colossal consideration as is classified as a measure that establishes the author, institutions and country (Borokhovich *et al.* 1995). Furthermore, cited articles have received more attention than less mentioned articles due to the influencing impact on the documents (Culnan, 1986). Again, h-index is modern techniques that combine both publications and citation counts under one framework (Hirsch, 2005). The bibliometric analysis uses various indicators in related identified scientific fields, which considered the database's pattern, such as most frequently cited articles, journals, production and cited (authors, institution and countries), and keywords through tables.

## 4 Findings and Results

## *4.1 Publication trends*

The study is base on 348 articles written on agency theory relating to financial and management accounting published between 1999-2019. **Figure 1** shows years of publications over the period. The graph shows that the selected sample of publication activities around the scientific field started in 2010 with 17 papers, and have increased substantially from 2015 onwards. The possible reason for constant growth for this fact might be that accounting information may serve as a source of mediation between shareholders and managers in communicating a company's

position and performance. Managers provide a company's accounting situation to shareholders might be relevant and useful for decision-making through a sound accounting information system, which has become one significant aspect of today's agency contract relationship and today's research activities.



**Figure 1 No.** Annual Publications in Agency Theory in Accounting

*4.2 Leading journals*

To identify which academic journal most frequently and impactful is essential in advancing scientific discipline in a related topic, as to reflect priority areas for scholars discussion and researchers, and to help other practitioners selecting which journal is appropriate for contributions to their manuscript. The 10 chosen journals including Critical Perspective on Accounting(CPA), Accounting Auditing and Accountability(AAA), Academy of Accounting and Finance Studies (AAFS), Australian Accounting Business and Finance(AABF), Corporate Ownership and Control(COC), European Accounting Review (EAR), Auditing (A), Accounting and Business Review(ABR), Accounting Organisation and Society (EAR), Accounting Research(AR), as presented in

Table **1**. The top 10 selected sample articles are CPA and AAA are the most published journal with 10 and 9, respectively. At the same time, AAFS, AABF and COC journals

63

with eight publications each, are ranked 3rd. The rest of the journals, all but ABR, AOS and AR, are rank 5 the least of the sample.

On the citation, AAA journal,  with 9 papers, has the highest citation count of 1001 with h-index of 4.37. Subsequently, AR, AOS, CPA, A, have a substantial citation count of 671, 503, 472, 320 respectively from the sample. Furthermore, other journals EAR, ABR, AAFS and AABF also have notable citation counts. The journal of Auditing review among the least production is the flagship journal which has the highest h-index. It means that researchers consider 4 publications as influential in the scientific field.

**Table 1. 10** Leading Journals in Agency Theory in Accounting

| Name | Abbreviations | Total Production | Total Citation | H-Index |
|---|---|---|---|---|
| Critical Perspective on Accounting | CPA | 10 | 472 | 4.10 |
| Accounting, Auditing And Accountability Journal | AAA | 9 | 1001 | 4.37 |
| Academy of Accounting And Financial Studies Journal | AAFS | 8 | 177 | 0.79 |
| Australasian Accounting Business And Finance Journal | AABF | 8 | 113 | 1.10 |
| Corporate Ownership And Control | COC | 8 | 61 | 0.11 |
| European Accounting Review | EAR | 7 | 294 | 2.83 |
| Auditing | A | 5 | 320 | 2.86 |
| Accounting And Business Research | ABR | 4 | 274 | 2.54 |
| Accounting Organisation and Society | AOS | 4 | 503 | 4.26 |
| Accounting Research | AR | 4 | 671 | 6.71 |

### 4.3 Cited articles

**Table 1** the most 10 cited articles on the scientific topic, group of authors, title, year, journal and citation. The most highly cited paper was published in 2012 by Cho, Freeman and Patten. The article has received 92 citations and disseminated from the leading journal's AAA in

Table **1**. The authors see AAA journal as influential in posting a scientific paper. In the same publication year, three are other documents received a considerable influence with 28, 27 and 20 as identified. The 2012 articles have dominated about 33.33%. Similarly of the top 10 cited articles, 4 publications are in the year 2011 with citation 52, 46, 43 and 12 from respective journals, take a share of 33.3%. For all the 10 top-

cited articles, in the years of 2014, the citation is 38, and 2015 is 34, then, the percentage received is 8.33% each.

**Table 2 10 cited articles**

| Author/s | Title | Year | Journal | Citation |
|---|---|---|---|---|
| Cho C.H., Freedman M., Patten D.M. | Corporate disclosure of environmental capital expenditures: A test of alternative theories | 2012 | Accounting, Auditing and Accountability Journal | 92 |
| Whittington R. | The Practice Turn in organisation research: Towards a disciplined transdisciplinarity | 2011 | Accounting, Organisations and Society | 52 |
| Kilfoyle, E., Richardson A.J. | Agency and structure in budgeting: Thesis, antithesis and synthesis | 2011 | Critical Perspectives on Accounting | 46 |
| Østergren K., Stensaker I. | Management control without budgets: A field study of 'Beyond Budgeting' in practice | 2011 | European Accounting Review | 43 |
| Vosselman E | The performativity thesis and its critics: Towards a relational ontology of management accounting | 2014 | Accounting and Business Research | 38 |
| Trotman A.J., Trotman K.T. | Internal audit's role in GHG emissions and energy reporting: Evidence from audit committees, senior accountants, and internal auditors | 2015 | Auditing | 34 |
| Niemi L., Kinnunen J., Ojala H., Yroberg P | Drivers of voluntary audit in Finland: To be or not to be audited? | 2012 | Accounting and Business Research | 28 |
| Quagli A., Avallone F. | Fair value or cost model? Drivers of choice for IAS 40 in the real estate industry | 2012 | European Accounting Review | 27 |
| Hyvönen T., Järvinen J., Oulasvirta L., Pellinen J. | Contracting out municipal accounting: The role of institutional entrepreneurship | 2012 | Accounting, Auditing & Accountability Journal | 20 |
| Moore D.R.J. | Structuration theory: The contribution of Norman Macintosh and its application to emissions trading | 2011 | Critical Perspectives on Accounting | 12 |

## 4.4 Productive and cite authors

Table **3**, presents the 10 most global productive and cited authors of the scientific field statistics, and a considerable h-index over the period. It has also included the country and the institution of origin. Considering the total number of papers publication, Modell has the highest publications of articles in the related area. Craig

65

and Jack follow the result have 3 publications each on the field. Further observation revealed that the remaining authors have 2 publications each over the period.

On the issue Of citation, it is surprising that Craig with 3 publications is the highest cited author and highly indexed factor of 28. Modell with 1827 citation count follows this. Not far behind him are Hussainey and Freeman with citation counts of 1393 and 1283 respectively. Moreover, Dhliwayo with 2 papers from South Africa had no citation and h-index among the top 10 influential authors.

Out of 10 most productive authors, 4 are from United Kingdom institutions, Modell, Craig, Jack and Hussainey. Furthermore, three influential authors affiliated with United Kingdom institutions. The top remaining authors are from other institutions and countries in the USA, Canada, Malaysia and South Africa. This result is exciting because an author from South Africa with 2 papers did not receive either citation or h-index.

**Table 3 10 Most Productive and Influential Authors**

| Authors | Total Productions | Total Citation | H-Index | Country | Institution |
|---|---|---|---|---|---|
| **Modell, S** | 4 | 1827 | 23 | United Kingdom | Alliance Manchester Business School |
| **Craig, R** | 3 | 2093 | 28 | United Kingdom | Durham University of Business School |
| **Jack, L** | 3 | 312 | 9 | United Kingdom | University of Portsmouth |
| **Brivot, M** | 2 | 204 | 8 | Canada | Universite Laval |
| **Bussin, M.H. R** | 2 | 46 | 4 | South Africa | University of Johannesburg |
| **Dhliwayo, D. V** | 2 | 0 | 0 | South Africa | University of Pretoria |
| **Freeman, M** | 2 | 1283 | 15 | United State of America | Townson University |
| **Himick, D** | 2 | 86 | 5 | Canada | University of Ottowa |
| **Hussainey, K** | 2 | 1392 | 21 | United Kingdom | University of Portsmouth |
| **Kallamu, B. S** | 2 | 10 | 1 | Malaysia | University of Putra Malaysia |

## 4.5 Authors cited journals.

Table **4** shows lists of 10 most productive and influential authors who have contributed papers in the top selected journals on procedure agency theory integrating into accounting. The 24 articles from the 10 journals have publications

and citations in the related scientific field. The highest productive author has 4 papers from AAA, EAR, AOS journals and ranked the second-highest cited author. However, the highest cited author Craig with 2093 (3apapers), did not have an article among the top listed journals. Besides this, Bussin and Dhliwayo also have no publications from any of the leading journals. It means that it could be acceptable among other journals. The remaining authors have quite a few published papers in the journals and citations, which is quite impressive.

Concerning authors contributed to the top journals, AAA journal contains the highest publications of 4 and most cited

Table **1**, followed by EAR and AOS journals with 3 papers each from the top authors. The next relevant journals (COC and A R) have 1 article each with their authors included in 10 top productive and cited authors. Nevertheless, remaining five journals (CPA, AAFS, AABF, A, ABR) have no article from leading authors, meanwhile are included in the top sampled journals in

Table **1**.

**Table 4 Productive Authors in the 10 Journals of Agency theory in Accounting**

| | TP | TC | CPA | AAA | AAFS | AABF | COC | EAR | A | ABR | AOS | AR |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Modell, S | 4 | 1827 | | 1 | | - | | 1 | | | 1 | |
| Craig, R | 3 | 2093 | - | - | - | - | - | - | - | - | - | - |
| Jack, L | 3 | 312 | - | 2 | - | - | - | - | - | - | - | - |
| Brivot, M | 2 | 204 | - | - | - | - | - | 1 | - | - | 1 | - |
| Bussin, M.H.R | 2 | 46 | - | - | - | - | - | - | - | - | - | - |
| DhliwayO, D.V | 2 | 0 | - | - | - | - | - | - | - | - | - | - |
| Freeman, M | 2 | 1283 | - | 1 | - | - | - | - | - | - | - | - |
| Himick, D | 2 | 86 | - | - | - | - | - | 1 | - | - | 1 | - |
| Hussainey, K | 2 | 1392 | - | - | - | - | - | - | - | - | - | 1 |
| Kallamu, B.S | 2 | 10 | - | -- | - | - | 1 | - | - | - | - | - |

## 4.6 Productive and influential institutions

Another distinct feature of the bibliometric analysis is the institutions that are interested in the publications of the scientific field of research. **Table 5** presents a list of the institution contributions in the number of publications, citations, and highly index. The study shows that the University of Manchester is the most productive institution. It has 10 publications, ranked second highest cited institution with 97, and with the highest h-index of 6. After this, Alliance Manchester Business School is

the second most productive institution in terms of papers on the scientific field, of his 9 publications ranked third cited institution with 93 counts, at the same time with 6 h-index among the top. The University of Texas, the only institution from the USA, has 6 publications. Of these publications, it has received 109 citations, the highest among the top 10. It means the impacting factor from the articles is higher. The rest of the institutions next to it are the University of Utara, Victoria University of Melbourne, University de Minho and the University of Portsmouth have 5 papers each with corresponding citation and h-index. Tampereen University is among the least published institutions with no citation count, and it is probably new in Scopus. All the institutions are from the United Kingdom and have 4 out of the 10 top institutions.

**Table 5 The Most Productive and Influential Institutions**

| Institution | Country | Total Production | Total Count | H-Index |
|---|---|---|---|---|
| University of Manchester | United Kingdom | 10 | 97 | 6 |
| Alliance Manchester Business School | United Kingdom | 9 | 93 | 6 |
| University of Texas Austin | United State of America | 6 | 109 | 3 |
| University Utara Malaysia | Malaysia | 5 | 6 | 2 |
| Victoria University of Melbourne | Australia | 5 | 52 | 5 |
| Universidade do Minho | Portugal | 5 | 16 | 2 |
| University of Portsmouth | United Kingdom | 5 | 24 | 2 |
| University Sains Malaysia | Malaysia | 4 | 9 | 2 |
| University of Essex | United Kingdom | 4 | 47 | 3 |
| Tampereen Yliopisto | Finland | 4 | 0 | 3 |

## 4.7 Country analysis

Table 6 presents a total of 10 countries of origin where the publications are studied. The objective is to see the volume of production and the most influential country because this reflects the importance of a country's contribution to the subjective matter. The United States of America(USA), a country with the most product, with 98 publications, ranked highly impactful state with 1422 citation counts and indexed 19 in the scientific field, far above the United Kingdom(U K) about 46% in terms of

publications, 106% in citations and 46.15% in h-index. Australia (38), China(25) and Malaysia(21) ranked chronological as the next in productivity. Other top countries among the list Canada and Germany, have 12 papers from each country, representing sixth in the line of publication. Same as in the case of France and Indonesia with 11 articles found as countries with substantial publications. Italy ranked the least publisher on the scientific field and is among the rest countries with most impactful in the area of study.

**Table 6** **The 10 Most Productive Country**

| Country | Total Production | Total Citation | H-Index |
|---|---|---|---|
| **United State of America** | 98 | 1422 | 19 |
| **United Kingdom** | 67 | 689 | 13 |
| **Australia** | 38 | 343 | 11 |
| **China** | 25 | 111 | 6 |
| **Malaysia** | 21 | 53 | 4 |
| **Canada** | 12 | 97 | 5 |
| **Germany** | 12 | 69 | 6 |
| **France** | 11 | 212 | 6 |
| **Indonesia** | 11 | 6 | 2 |
| **Italy** | 10 | 145 | 6 |

## 4.8 Keywords used in agency theory in accounting

Keywords are part of the choice of the authors to represent the content of the article, be as general and ordinary as possible (Curty and Boccato, 2005; Comerio and Strozzi, 2019)**.** They provide additional steps to assess the information flow and trace any scientific study(Madani and Weber, 2016). The top keywords found from Scopus, which have frequently used in 348 articles in accounting agency theory, are represented. It shows the number of occurrences authors have utilised in the reference. The Agency theory has significant influence in accounting as many have frequently related to the scientific field. Accounting also is trendy among the top keywords used comparing with the other keywords, except corporate governance. Agency theory and accounting have impactful influence over a decision as they ranked among the top three in the table. Besides, they have become a novel focus with the rising becoming famous for researchers as impressive is increasing in the publication from 2015-2019.

**Table 7** **10 Keywords Used**

| Keywords | Frequency | Ranking |
|---|---|---|
| Agency Theory | 67 | 1 |
| Corporate Governance | 49 | 2 |
| Accounting | 19 | 3 |
| Agency | 17 | 4 |
| Earning Management | 15 | 5 |
| Firm Performance | 15 | 5 |
| Accountability | 11 | 6 |
| Human | 11 | 6 |
| Institutionary Theory | 11 | 6 |
| Ownership Structure | 9 | 7 |

## 5.0 Conclusion

In conclusions, the bibliometric analysis shows various kind of academic research that has gone on viral in the subject area of agency theory in accounting. The study outcomes were from the Scopus database as it considered as the most productive repository. The development from authors, countries, institutions, and publications from articles received more attention. It anticipated that future research should continue to bring more uncovered areas and conducted research in other disciplines to get more empirical evidence for academic development. The study contributes to academics and practitioners about the newer method for researching using bibliometric analysis tools.

The growth of the scientific field is receiving consideration substantially. The selected sample period started at 17 articles in 2010, and after that, it has consistently maintained publication moved to an average of 35 per year during the period 2010-2019. The result implies that researchers paid attention to the agency theory's impact in accounting though there were ups and downs(Figure 1).

The development sees that out of 10 most top journals that published an article on the scientific field, CPA is the highest production of 10 publications. Besides, AAA is the leading influential journal among the top 10 in terms of citation count of 1001. It has interdisciplinary combinations, and AR gets the best h-index of 6.71 in the research field. These attributes further strengthen each journal's point to receive global recognition from academia (see

**Table** **1**). Other journals received lower than 5 publications, even though receive substantial citation counts above 270-671.

Focusing on the highly cited article, the highest received 92 citations for the most part referenced paper compare with others from AAA journal in 2012. The next leading articles that received high cited count happen in 2011 range from 43-52. Others received related citations, as shown in (

**Table 2**).

In terms of the most productive and influential author, the analysis highlighted Modell with 4 papers as the most productive author, followed by 2 authors Craig and Jack with 3 articles each of the scientific field and 7 least authors with 2 papers each. Besides, Craig is the highest cited and highly indexed author among them,

**Table** **3**. Furthermore, 4 of the top authors are from the United Kingdom and have published most highly cited papers in the field.

Another contributory to the research is top authors in the top identify journals. Modell has 3 publications, one each in the following journals AAA, EAR and AOS. Craig, highly cited and h-index author and other 3 authors, do not have any publications in the top journal in

**Table 4**.

The United Kingdom institutions have shown majority position in the production in the field understanding, the agency theory in accounting. The institutions apart from Texas Austin 109 cited counts is the second most influential ones shown in Table 5. They have 28 publications out of 57 papers in the top 10 institutions. Furthermore, it has many of the leading authors.

The United States is the most dominant country in the scientific field in terms of article publication, highly influential nation, and very well indexed on the list. The position followed by the United Kingdom with 67 papers ranked second highest cited. Its impacting factor is remarkable compared to other countries whose contributions in the scientific field are not near them (Table 6)**.**

Furthermore, another feature contribution to the scientific field is keywords. It emerges that 69 agency theory has been frequently cited and 19 times from accounting. It means that the subject is directly related to the subject area (.

**Table 7**).

This paper's main results are useful for obtaining, based on bibliometric data, a general overview of the state -of- the -art research regarding agency theory relation in accounting. It also directs future research efforts on innovation adoption by offering a broad understanding of using extant literature review and bibliometric analysis tools on the current trend of research in different contexts and disciplines. However, there are limitations which can offer for future directions. One of the study's obstacles was choosing the database. First of all, relevant and quality agency theory in accounting research performed beyond indexed journals that not included in the study. Second, the sample selection is limited to unique search keywords, and entries, although attempts to ensure that all related publications are selected. Thirdly, the citations represent the influence level of a paper, author or journal at a particular time. The citation count can increase considerably over time, representing different numbers and ranks. Finally, this study contains a general overview that can help clarify the scientific field of agency theory in accounting. Still, several other topics need to be taken into account to get a full picture of the art state.

Future research into areas using visualisation of co-authorship network of countries and institutions, co-citation networks of publications and journals on the scientific field proposed. Moreover, studies can expand the scope by considering other types of review papers and documents to incorporate more detailed information for the study field. Besides, similar studies at future stages should continue to reveal the literature review trend on agency theory in accounting and track its continuing credibility and development.

## References

Anthony, R. N. (1965). *Planning and control systems: A framework for analysis* (p. 17). Harvard Graduate School of Business.

Baiman, S., Netessine, S., & Saouma, R. (2010). Informativeness, incentive compensation, and the choice of inventory buffer. *Accounting Review*, *85*(6), 1839–1860.

Beest, F., Braam, G., & Boelens, S. (2009). *Quality of Financial Reporting: Measuring qualitative characteristics. NICE.* (No. 09–108).

Berger, A. N., & Udell, G. . (1988). "The Economics of Small Business Finance: The Roles of Private Equity and Debt Markets in the Financial Growth Cycle." *Journal of Banking and Finance*, *22*, 613–673.

Berger, A. N., & Udell, G. . (2006). "A More Complete Conceptual Framework for SME Finance." *Journal of Banking and Finance*, *30*(11), 2945–2966.

Berle, A., & Means, G. (1932). *The modern corporation and property*. Macmillan.

Borokhovich, K. A., Bricker, R. J., Brunarski, K. R., Simkins, & B.J. (1995). Finance Research Productivity and Influence. *Journal of Finance*, *50*(5), 1691–717.

Bouckova, M. (2015). Management Accounting and Agency Theory. In D. Prochazka (Ed.), *16th Annual Conference on Finance and Accounting (ACFA)* (pp. 5–13).

Brink, A. G., Hobson, J. L., & Stevens, D. E. (2017). The effect of high power financial incentives on excessive risk-taking behavior: An experimental examination. *Journal of Management Accounting Research*, *29*(1), 13–29.

Brown, L. D. (1996). "Influential accounting articles, individuals, Ph.D. granting institutions and faculties: A citational analysis." *Accounting, Organisations & Society*, *21*(7/8), 723–754.

Brown, P., Beekes, W., & Verhoeven, P. (2011). Corporate governance, accounting, and finance: A review. *Accounting and Finance*, *51*, 96–173.

Cavélius, F. (2011). Opening the ―black box‖: How internal reporting systems contribute to the quality of Abstract financial disclosure. *Journal of Applied Accounting*, *12*(3), 187–211.

Chi, S.-K. (1989). *Ethics and Agency Theory*. An Arbor, MI.

73

Cobo, M. J., Lopez-Herrera, A. G., Herrera-Viedma, E., & Herrera, F. (2011). Science mapping software tools: Review, analysis, and cooperative study among tools. *Journal of the American Society for Information Science and Technology*, *62*(7), 1382–1402.

Comerio, N., & Strozzi, F. (2019). Tourism and its economic impact: A literature review using bibliometric tools. *Tourism Economic*, *25*(1), 109–131.

Culnan, M. J. (1986). "The Intellectual Development of Management Informations Systems, 1972-1982: A Co-citation Analysis." *Management Science*, *32*(2), 156–172.

Cunha, P. ., Toigo, L., & Picoli, M. . (2016). Scientific production on the Audit Committee: A Bioblimtric and Sociometric Analysis of International Journals. *DIGITAL LIBRARY OF JOURNALS*, *8*(1), 26–46. https://doi.org/http://dx.doi.org/10.5380/rcc.v8i1.35902

Curty, M. G., & Boccato, V. R. C. (2005). O artigo científico como forma de comunicação do conhecimento na área de Ciência da Informação. *Perspectivas Em Ciência Da Informação (Impresso)*, *10*(1), 94–107.

Duan, W. (2014). Research on the budgetary slack behavior to the construction machinery and equipment manufacturing enterprise in China. *Advanced Materials Research*, *915–916*, 1504–1508.

Eisenhardt, K. M. (1989a). Agency theory: An assessment and review. *Academy of Management Review*, *14*(1), 57–74.

Eisenhardt, K. M. (1989b). Agency theory: An assessment and review. *Review, Academy of Management*, *14*(1), 57–74.

Harrs, M., & Raviv, A. (1978). Some results on incentive contracts with applications to education and employment, health insurance, and law enforceement. *American Economic Review*, *68*(1), 20–30.

Hiebl, M. R. . (2015). Agency and stewardship attitudes of chief financial officers in private companies. *Qualitative Research in Financial Markets*, *7*(1–2), 4–23.

Hirsch, J. . (2005). An Index to quantify an Individual's scientific research output. *Proceedings of the National Academy of Sciences of the United States of America*, *102*(46), 16569–16572.

Huang, R., Marquardt, C. A., & Zhang, B. (2014). Why do managers avoid EPS dilution? Evidence from debt-equity choice. *Review of Accounting Studies*, *19*(2), 877–912.

Ijiri, Y., Levy, F. K., & Lyon, R. . (1963). A linear programming model for budgeting and financial planning. *Journal of Accounting Research*, *1*(2), 198–212.

Jeggers, M. (2008). *Managerial economics of non-profit organisations (Routledge studies in the management of voluntary and non-profit organisations*. Routledge.

Jensen, M. C., & Meckling, W. . (1976). Theory of the firm: managerial behavior, agency costs and ownership structure. *Journal of Financial Economics*, *3*(4), 305–360.

Kluvers, R., & Tippett, J. (2011). An exploration of stewardship theory in a Not-for-Profit organisation. *Accounting Forum*, *35*(4), 275–284.

Lambert, A. R. (2007). Agency theory and Management Accounting. *Handbook of Management Accounting Research*, *1*, 247–267.

Madani, F., & Weber, C. (2016). The evolution of patent mining: Applying bibliometrics analysis and keyword network analysi. *World Patent Information*, *46*, 32–48.

Maksimov, V., Wang, S. L., & Luo, Y. (2017). Institutional imprinting, entrepreneurial agency, and private firm innovation in transition economies. *Journal of World Business*, *52*(6), 854–865.

Martin, G. ., Thomas, W. ., & Wieland, M. . (2016). S&P 500 Membership and Managers' Supply of Conservative Financial Reports. *Journal of Business Finance and Accounting*, *43*(5–6), 543–571.

Merigo, J. M., Gil-Lafuente, A. M., & Yager, R. R. (2015). An Overview of Fuzzy Research with Bibliometric Indicators. *Applied Soft Computing*, *27*, 420–433.

Mitnick, B. (1975). The theory of agency: The policing, "paradox" and regulatory behaviour. *Public Choice*, *24*(1), 27–42.

Podsakoff, P. ., MacKenzie, S. ., Podsakoff, N. ., & Bachrach, D. . (2009). 'Scholarly Influence in the Field of Management: A Bibliometric Analysis of the Determinants of University and Author Impact in the Management Literature in the Past Quarter Century.' *Journal of Management*, *34*(4), 641–720.

Pritchard, A. (1969). Statistical bibliography or bibliometrics? *Journal of Documentation*, *25*(4), 348–349.

Qu, S. ., Ding, S., & Lukasewich, S. . (2009). 'Research the American Way: The Role of US Elites in Disseminating and Legitimising Canadian Academic Accounting Research',. *European Accounting Review, 18*(3), 515–69.

Ramos-Rodrıguez, A. R., & Ruız-Navarro, J. (2004). Changes in the intellectual structure of strategic management research: A bibliometric study of the Strategic Management Journal, 1980–2000. *Strategic Management Journal*, *25*(10), 981–1004.

Ross, S. (1973). The economic theory of agency: The principal's problem. *American Economic Review*, *63*(2), 134–139.

Schoenfeld, J. (2020). Contracts Between Firms and Shareholders. *Journal of Accounting Research*, *58*(2), 383–427.

Schroeck, G. (2002). *Risk Management and Value Creation in Financial Institutions* (G. Schroeck (ed.)). Wiley.

Scott, W. R. (2015). *Financial Accounting Theory* (7th ed.). Ontario Pearson Canada.

Steinberg, R. (2010). *Principal–agent theory and nonprofit accountability*. Cambridge University Press.

Stiglitz, J. E., & Weiss, A. (1981). "Credit Rationing in Markets with Imperfect Information." *The American Economic Review*, *71*(3), 393–410.

Tang, M., Liao, H., & Su, S. . (2018). A Bibliometric overview and visualisation of internaltional Journal of Fuzzy Systems between 2007-2017. *International Journal of Fuzzy System*, *20*(5), 1403–1422.

Zulkarnaen, W., Bagianto, A., Sabarb, & Heriansyah, D. (2020). Management accounting as an instrument of financial fraud mitigation. *Nternational Journal of Psychosocial Rehabilitation*, *24*(3), 2471–2491.

Chapter 3

# Using Agency Theory to Model Cooperative Public Purchasing

*Cliff McCue and Eric Prier*

## INTRODUCTION

The operational linkages between government organizations, their purchasers and their suppliers are now recognized as important contributors to the success of government policy and decision-making. Although cooperative purchasing has been a topic of study for many years (Anderson & Macie, 1996; Goodwyn, 1976; Johnson, 1983; Knapp, 1969; Miller, 1937; Saloutos & Hicks, 1951; Taylor, 1953), researchers have only recently revisited issues related to cooperative public purchasing (CPP) in search of more clarification with respect to its theoretical underpinnings (Aylesworth, 2003; Tulip, 1999; Wooten, 2003). Perhaps due to little theoretical direction and few standards to guide practice, there seems to be a lack of conceptual coherence within the cooperative purchasing literature to inform us concisely about what comprises cooperative procurement and its implications for public purchasing. Indeed, John Ramsay and Nigel Caldwell (2004) make a strong case that metaphors so often used can lead to misunderstanding the nature of interesting phenomenon. It is no different in public purchasing, as slight misconceptions about institutional goals and to whom one is accountable may in fact have significant organizational consequences. To help remedy this situation, this chapter provides a long-needed general framework in utilizing agency theory to analyze, define, and theoretically model cooperative public purchasing.

A theory of CPP is needed for at least two reasons. First, a theory in this area can help all stakeholders in public procurement better understand the role they play in providing incentives for utilizing cooperatives in purchasing decisions. For example, the limited amount of extant research indicates that the term itself is conceptually muddled. Without a systematic theory to offer guiding principles for the phenomena called cooperative purchasing, the imprecise cognitive

images bandied about by practitioners and academics are merely that ambiguous notions of purchasing mechanisms that appear to be related in some way. Thus, consistent with the claims of academics that all theoretical perspectives have value (O'Toole, 1995; Terry, 1999), development of a theory of cooperative purchasing is useful if one wants to explain, predict and understand behavior concerning the intent, purpose, and actual use of cooperatives in procurement.

A second reason why theory is needed in this area relates to the first point: without a unified model, observers and practitioners alike must remain content with various depictions that appear to be cooperative purchasing. In fact, the small amount of research in this area is by its nature, merely informative. Without axiomatic and generalizable principles, practitioners, policymakers, and academia must make prescriptive recommendations without understanding the numerous potential consequences of engaging in CPP, or whether one model of CPP is better than another. This chapter is a first attempt to justify the use of specific conceptual terms which can structure these long-needed organizing principles, while providing direction to practitioners.

Given what is known about cooperative public purchasing at this critical point, agency theory holds considerable promise in connecting empirical observation to a generalizable theory of cooperative purchasing. Besides contributing new ways to approaching old problems, agency theory can help explain the purchasing incentives of individual purchasers by modeling their underlying motivations and clarifying the needs and goals of the stakeholders who support the cooperative purchasing process. In turn, this knowledge may be helpful in offering fundamental guidance to organizations that wish to transition from operational to strategic purchasing.

The current chapter is broken down into three sections. The first section discusses agency theory and its usefulness in modeling CPP, and briefly argues the need for conceptual clarity. The second segment explicates specific terms, definitions, and concepts that will be used to build three models of cooperative public purchasing within the context of agency theory. The third and final part discusses the strengths and weaknesses of the theory before laying out a research agenda based upon the models offered here.

## AGENCY THEORY

According to Jensen and Meckling (1976, p. 308), an agency relationship is "a contract under which one or more persons (principals) engages another person (the agent) to perform some service on their behalf which involves delegating some decision-making authority to the agent."[1] When executing the tasks within the principal-agent relationship, the agent must choose actions that have consequences for both the principal and the agent.[2] Since these outcomes can be either negative or positive for each of the actors, the chosen action of the agent affects the welfare of both. The principal-agent relationship is often forged because the agent possesses a greater abundance of the needed skills, abilities, and/or time to perform the desired activities. Inevitably, however, there are several problems for the principal in governing the relationship with the agent, the first of which involves choosing an appropriate agent.

Consistent with the tenets of agency theory, the view adopted here assumes that agents, purchasing officials, are rational, self-interested utility maximizers. However, it is not assumed that these agents behave selfishly and do so with guile. In other words, slightly contrary to Williamson's (1985) transaction cost economics framework, although it is assumed that people are opportunistic in the sense that they may shirk in a self-interested manner by trying to minimize effort if it fulfills their needs, it is not assumed that they will willingly misrepresent or lie about that effort. More to the point, it is merely assumed that the principal and agent do not share the same levels of information, and as such, the agent can opportunistically take advantage of the situation, sometimes to the detriment of the principal. This latter situation is known as moral hazard and is often the result of asymmetric information.[3]

### Asymmetric Information

Agency theory in economics has long been concerned with the issues of control that arise as a result of information asymmetries between agents delegated to maximize the welfare of the principals who contracted with them (see especially Ross [1973]; Jensen and Meckling [1976]). In general, all principal-agent relationships are plagued by uncertainty uncertainty not only in the level of an agent's knowledge, skills and abilities, but also in both the way the agent's action gets transformed into the output and whether or not the agent is acting in the

principal's best interest. This uncertainty is the result of the advantageous differential in knowledge held by the agent about his or her own actions in serving the principal. This difference is information asymmetry, and it is a third problem for the principal in governing the relationship with the agent.

Although under normal circumstances both the principal and the agent can observe the outcome, it is often the case that the principal cannot or does not observe the agent's specific action, effort, or capacity to perform all of which are supposed to obtain the outcome favored by the principal.[4] However, one must be cognizant that the agent not only observes her own action, but also may have knowledge not possessed by the principal about other factors that lead to the outcome. This information asymmetry describes the inability of the principal to properly assess the extent to which the agent chooses an action that coincides with the principal's best interests. As such, there can be little doubt that asymmetric information permeates the principal-agent relationship. For example, consider the case where a cooperative is used to purchase police vehicles. During the vendor selection process, purchasers may become aware of information that could potentially bias who would be selected in their own organization, such as satisfying local preferences, but given the nature of the cooperative, the agents do not divulge this information to their principals. The principals, on the other hand, may wish that a particular vendor be selected, so they reject the selection offered by the cooperative.

Not only do the actions or inactions of both the agent and principal influence the outcome, but also there are random factors, beyond the control of either the principal or the agent (which influence the outcome). Moreover, there are costs borne by the agent in performing the action, and by the principal in providing compensation in addition to the costs of monitoring the behavior of the agent. As such, these tools of agency theory are an appropriate lens by which to model cooperative public purchasing for at least three reasons. First, the nomenclature developed here can apply to both public and private cooperatives in a way that makes comparison easier between organizations in the two sectors, especially in identifying motivational similarities and differences.

Second but just as important, it is easy to see how there is a chain of agency relationships in cooperatives that may impact the nature of purchasing.[5] For example, procurement officials typically are employed

by each of the organizations who wish to cooperate, even while the cooperative association itself may or may not have a separate purchaser for the collective enterprise. Depending upon the rules, procedures, and by-laws of the cooperative agreements, there may be differential incentives for purchasers at various levels to utilize cooperatives that may not be apparent without the aid of agency theory.

Since CPP can be thought of as a chain of agency relationships similar to the contractual relations found within the economic firm, valuable questions arise as to the best way to organize the stakeholder relationships in public procurement. Ronald Coase (1937) was the first to reformulate the notion of the firm in orthodox economic theory from that conceived as a "black box" that transforms inputs (resources) into outputs (production). Instead, he conceived it as the neoclassical economics perspective of a system of relationships which directs production. This implies that a firm is more efficient at aligning resources with outputs than is the market. As Harold Demsetz (1983, p. 377) observes, "it is a mistake to confuse the firm of economic theory with its real world namesake. The chief mission of neoclassical economics is to understand how the price system coordinates the use of resources, not the inner workings of real firms." Similar to Coasian economics, procurement can be arranged through the market and regulated by the price mechanism with all of its attendant hidden costs to the procurement official, or the exchange transactions of procurement can be vertically integrated and ordered through the firm in a hierarchy where purchasing is integrated with the needs for the same products by other principals (and as we shall see, their agents).

This theoretical difference between market and hierarchy is not completely esoteric because the issues surrounding why exchanges take place in a market or under a firm is similar to discerning why procurement officials choose to cooperatively buy. Clearly it is not costless to find a good cooperative to use, and it is important to understand what benefits accrue to procurement officials and their principals to explain why they take on the additional costs associated with utilizing CPP. Thus agency theory can expose the motivations of stakeholders in public procurement.

As mentioned previously, agency theory and its embedded theory of incentives generally assumes that actions and efforts are normally unverifiable, while outcomes are generally known and confirmable (Dixit, 2002, p. 713). In terms of CPP, the effort of the procurement official is

verified only when the outcome (the purchase) is obtained. However, it will be shown that the action, as opposed to the outcome, may not be readily distinguished by the stakeholders. Consider that although the procurement official might believe that the actual purchase is an "outcome," the purchase is merely considered an "action" from the viewpoint of the stakeholder for whom the purchase was made. In other words, the level of analysis is important in determining what behavior is an "action" as opposed to an "outcome."

The third reason why agency theory is a fruitful method for modeling cooperative public purchasing is that it helps to identify the various incentives of the stakeholders. By clarifying the opportunities and constraints they face, hope is engendered that efficiency, effectiveness, and accountability will be increased.

## LEVELS OF ANALYSIS

Models are simple approximations of a given phenomenon, and when examining and modeling public purchasing cooperatives, the level of analysis is extremely important in determining the conceptual attributes of interest. For example, cooperatives can be defined on at least four levels. At one level, the independent government entities engaged in cooperative purchasing would be the focus. These entities are termed public cooperative affiliates (PCA), and they comprise the members of the cooperative. At this organizational level, it is assumed that government entities require purchasing departments and organizational personnel to coordinate their purchasing activities in a way that is relevant to each PCA participating in the cooperative. At this level of analysis, the government itself is the principal that relies on both the purchasing departments and the purchasing officials who are the focus of the next level.

At another level, even the individual purchasers within those government agencies could be the effective unit of analysis. In this way of thinking, the public purchasers would be interdependent actors who are asymmetrically informed about the costs, benefits, and management of the cooperative enterprise, and it may be that they each have different motivations for utilizing cooperative purchasing. This suggests that purchasing officials might shirk their responsibilities to citizens, government, or their own purchasing department by diminishing effort or transferring effort to the cooperative. Moreover, some might be prohibited from exploiting CPP or directed to utilize a particular

cooperative organization for its purchasing needs, while others may have discretion in deciding if and when a cooperative is used. Still others may join a cooperative for some types of purchases and not for others.

At a more holistic level, the cooperative as a whole can be modeled as the unit of analysis. This would consist of the cooperative enterprise, the PCAs, and the organizational charter or the legal covenants governing cooperative public purchasing. These latter elements of the model are referred to as the cooperative public purchasing agreements (CPPAs), and at a minimum they should delineate four elements of the cooperative including: identifying who does the negotiating and buying (e.g., hereafter called the mechanism of purchasing); formalizing the organizational and institutional contours of the cooperative enterprise; specifying the dues paid by the PCAs for maintaining the CPPA; and outlining the relationship of the public cooperative affiliates to each other. Consequently, even though cooperatives are composed of bureaucratic organizations and individuals, this view sees the cooperative as a corporate body where the cooperative interacts intra-organizationally and with other entities in its environment. In other words, the entire cooperative integrates the agreements outlining the complete relationship of component members (PCAs) to the cooperative enterprise and its organization (if there is one).

In this attempt to model cooperative public purchasing at this level, what becomes important is not only understanding the interdependency among participating PCA members, but also determining CPPA organizational responsiveness, transparency, and alignment of goals with member PCAs by ascertaining the negotiation and purchasing procedures for the CPPA. Finally, one could model and define cooperative public purchasing in terms of the social system whereby citizens, businesses, governments, and all potential vendors and suppliers could be mapped into a supply line, network, or web that focuses on modeling how the cooperative interacts with these and other societal stakeholders. At this level, there are numerous layers of agency and common agency.

## THE MODELS

A simple principal-agent theory of cooperative public purchasing is a powerful tool to view cooperative purchasing arrangements because it can be used to study purchasing process outcomes, stakeholder

behavior, information dissemination, decision-making, and accountability in cooperative arrangements. According to the underlying theory of the models, the principal is a stakeholder that retains a person or organization to undertake a specific task and serve a particular functional role within cooperative public purchasing. In turn, the person or organization delegated to manage these responsibilities on behalf of the principal is the agent.
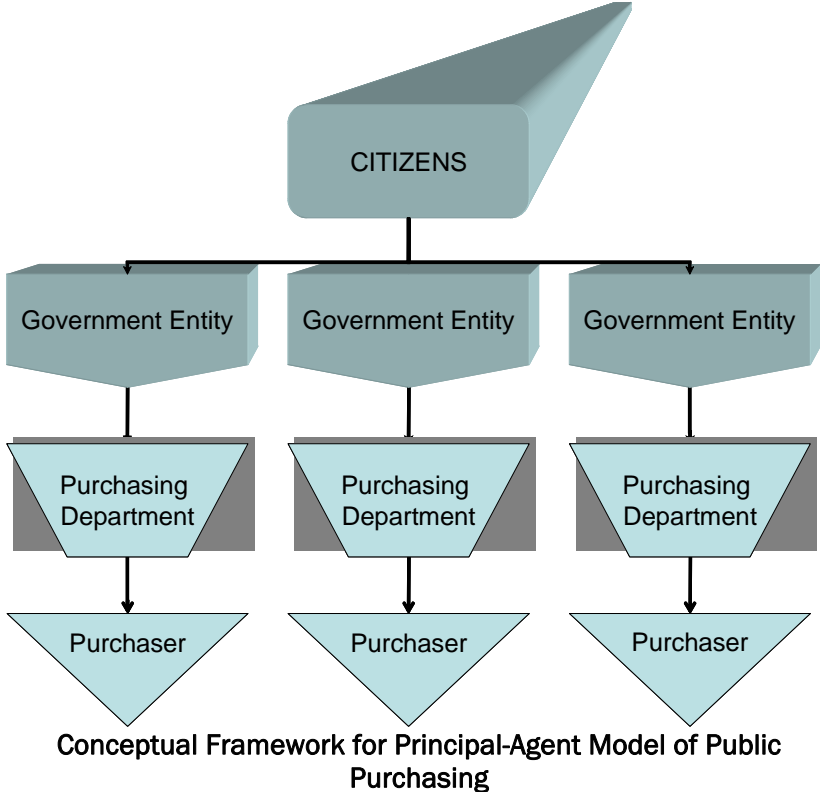
Although operationally, practitioners and theoreticians are most interested in the mechanism by which goods and services are purchased and the relationship of the affiliates (PCAs) to one another, there may be other considerations. For example, they might also want to know about how title passes from supplier to purchaser, the scope of purchases by the cooperative agreement (i.e., single-purpose or multi-purpose), the determination of the sharing of expenses, contracting issues such as the procedures for negotiating purchases, and questions about ownership of the cooperative (if there is ownership). However, for the sake of parsimony, two cooperative purchasing dimensions which appear to be basic elements of all cooperatives are discussed herein, and they are the mechanism of purchasing the actual goods and services, and the relationship of the affiliates (PCAs) to each other. If the mechanism of purchasing is located in an organ external of the cooperative itself, is it for-profit or non-profit? These are some issues to be explored.[6]

At its most elemental level, the model of public purchasing in representative democracies is depicted in Figure 1, suggesting a very simple model of a principal-agent hierarchy for purchasers in the public sector. It is assumed that government exists for the benefits of its citizens and thus is the agent of the people, the beginning and original principals. Moreover, it indicates that various government entities fulfill their own purchasing needs by utilizing their respective purchasing departments.[7] Thus, the departments become agents of the government entity, and the government entity is an agent of the citizens. However, there is a third level of agency depicted in this chain of agency, and it consists of the purchasing department employees (labeled purchaser) who become direct agents of their respective departments and indirect agents of both the government entities and of citizens.

One might quibble that this model of public purchasing may appear to be too reductionistic, but it is useful by suggesting at least three significant characteristics of public purchasing. First, it is obvious that even a denuded representation like that shown in Figure 1 reveals

several layers of agency that are not readily apparent from casual

**FIGURE 1**



**Conceptual Framework for Principal-Agent Model of Public Purchasing**

observation. However, it is obvious that there is a chain of agency where at any given node of agency relation, there are multiple principals for whom the agent has a fiducial responsibility. In turn, this suggests that purchasing decisions may be more complex than generally recognized in the literature. A second reason this theory is important is that it illustrates the ambiguous nature of common agency in the public sector. No matter if it is conceived as delegated or intrinsic, when public purchasing decisions are made, the purchasing agent takes on fiduciary responsibilities of multiple and perhaps conflicting principals.

In Figure 1, consider who and what comprise the group called citizens. If it is assumed that they are domestic providers, this group

would also include the vendors and suppliers themselves with whom the purchasers at the bottom are contracting. Given that there are social goals beyond mere economics that are thought to be important when procuring public goods and services, the murkiness with which the public purchasing role can be viewed is considerable. Take, for instance, a common situation where a local government has a preference policy to buy locally whenever feasible. When this happens, it is difficult for the public to gauge the objectives and success of the purchase, because buying locally may pressure procurement costs upward, which might mitigate the goal of lowest cost, best value, or other efficiency goals. As a result, issues of public control and accountability emerge, especially when government agencies are pursuing multiple missions and there is a fuzziness surrounding public objectives (Dewatripont, Jewitt & Tirole, 1999).

Undeterred by the problems of common agency, economists have modeled organized interests as the principal and government as the agent (Grossman & Helpman, 1996, p. 753), and it is often supposed that interest groups or private corporations can asymmetrically bias public policy in their direction (Becker, 1983; Faith, Leavens & Tollison, 1982; Lohmann, 1998). But just as citizens do not universally agree on many goals, various governments can also have competing goals with one another. For example, the state legislature may have the immediate goal of lower taxes while some localities under its jurisdiction may want additional goods and services, yet they may be restricted by state regulation from taxing to provide for them. Nonetheless, although the people living in these cities are citizens of those municipalities, they are also citizens of the state. Assuming that these city-dwellers want the increase in services, the problem of common agency for a purchasing official becomes one of deciding which goal is more important – an often non-obvious choice.

The theory further suggests that whether or not these entities are of the same government (e.g., agencies within the same government) or represent different governments, the figure leaves open the potential that government entities may or may not have similar goals. Consequently, the agency relationships modeled here suggest that governments, if not outright competitive, could be at cross-purposes so that there may be times when cooperative public purchasing is not mutually advantageous. An example is when the U.S. General Services Administration (GSA) will restrict the availability of goods and services from the supply schedule if, when left open to be used by other entities,

the use of the schedule results in lower supply or higher prices for the U.S. government. In other words, if the federal government is adversely affected by other entities procuring material through its supply schedule, the available supply schedule will be shrunk by the federal government to capture the savings under the schedule. In summary, Figure 1 makes clear that for any purchasing decision, the multiple layers of principals and agents make organizational responsiveness and maintaining transparency difficult.

Before offering some specific models of cooperative public purchasing, there are other issues suggested by the simple model in Figure 1, because there are fundamental questions concerning how one should model the economics associated with public procurement (Demsetz, 1971; Telser, 1969; also see Lloyd [2000] on symptoms and treatments of contracting pathologies). For example, although it may be a major factor in obtaining cooperative public purchasing agreements, attempting to control and reduce production costs may be less important than the demand schedules or policy preferences of citizens.

In other words, it is often unknown how marginal price costs associated with publicly procured goods and services might impact the amount of purchased materials, and hence, the costs paid by public purchasers for those goods and services obtained through cooperative agreements. It may be that public pressures for a particular course of action requiring large purchase orders may override the rational cost-benefit calculus of the decision. Yet due to the fog of agency layers, it is difficult to either reward or punish this type of behavior. Consequently, based upon how one models cooperative public purchasing, seemingly innocuous decisions may take on added import. Consider the case when bidders and buyers' cooperatives are competing. Since they must estimate demand because the true demand is often unknown, this can lead to economic inefficiencies. However when one models public purchasing as if market or individual demand is known with certainty, it is not difficult to arrive at an economically efficient solution. These two circumstances have direct consequences for cooperative public purchasing.

Again, take the case of the GSA. If so many entities are using the federal schedule, the implication is that there is purchasing certainty and thus economic efficiency. Yet if entities can opt out on individual purchases in an *a la carte* manner, inefficiencies are introduced due to demand uncertainty. What is more, it may be that the benefits of

entering into a cooperative agreement may be asymmetrically distributed across not only the PCAs themselves, but also across their constituents, and this might include some potential suppliers. Thus, this should be but one element of the contractual obligations of the PCAs, and it should be precisely outlined in the agreement.
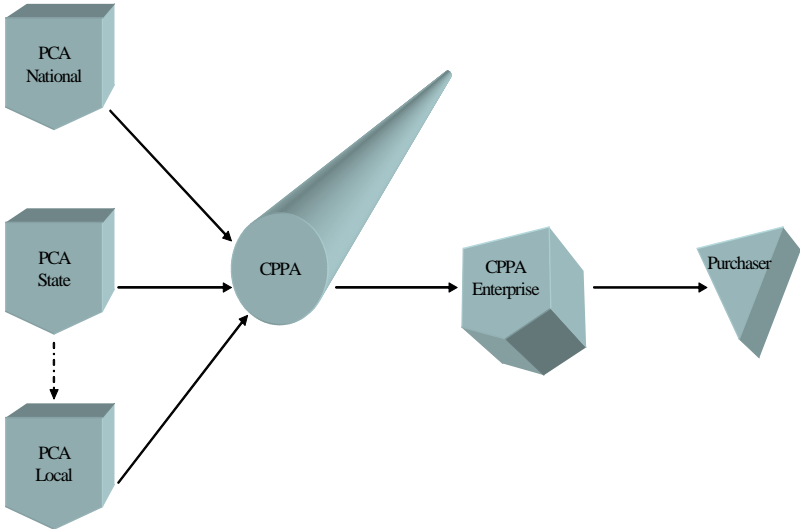
Indeed this realization is consistent with a U.S. Government Accounting Office (GAO, 1997) report which found that the effect of cooperative purchasing on industry providers and small businesses is likely to vary, and as the reinventing government movement advances, the problems of estimation are likely to become more complex. Moreover, geographic proximity of the cooperative public purchaser to the goods and service providers surely asymmetrically impacts the bottom line of the providers. Thus, the economies of transport costs are differentially enjoyed unless the cooperative agreement distributes the costs equally among the member affiliates. This helps to explain why a PCA may want to buy locally and in the GSA example mentioned above, why uncertainty and thus inefficiencies are introduced into cooperative public purchasing.

In terms of CPP, when one more clearly models and thus understands the principal-agent relationships associated with the PCAs, policymakers can more efficiently pursue clarified goals which are likely to result in significant savings of assets, resources, time and effort. Identifying rival providers, rival bidders, and potential non-rival partners can lead to cooperative agreements that make sense for all parties, while helping to fortify public purchasing. Thus, the models offered here will clarify relationships among institutions and individuals involved in cooperative public purchasing.

### Buyer Model

Figure 2 specifies a principal-agent model of cooperative public purchasing labeled the Buyer model. Under this system, the PCAs choose to promulgate a cooperative public purchasing agreement (labeled CPPA) which specifies that an administrative bureaucratic organ (labeled CPPA Enterprise) will be created to carry out the mandates of the CPPA. Charged with fulfilling this role, the CPPA enterprise itself hires individuals (labeled purchaser) to negotiate purchases and contracts for the membership. One example of the Buyer model is the Educational and Institutional Cooperative Service, Inc., which is a non-profit buying cooperative owned by more than 1,500 member-PCAs.[8]

FIGURE 2
Buyer Model of Cooperative Public Purchasing



There are several characteristics worthy of discussion in this simple Buyer model. First, even after excluding the citizens and the purchasing departments within the PCAs themselves, which adds additional layers of agency, the configuration readily shows the chain of agency just at the cooperative level. For example, the member-PCAs that coalesce around the CPPA utilize the enterprise as its agent, which in turn utilizes another agent (the purchaser). Thus, the purchaser becomes an indirect agent two steps removed from the original principals in the model! In this circumstance, issues of citizen accountability and control arise because the ultimate goal that is supposed to be served by the purchaser may be convoluted and recast as it moves through the chain. Under these conditions, the agent is the purchaser of the cooperative, and the principals are the various participants in the cooperative typically a public purchasing agent from one government.

A second quality of the Buyer pattern of public purchasing has to do with the relationship between PCA's. The dotted line represents the possibility of a vertical public purchasing agreement, which describes
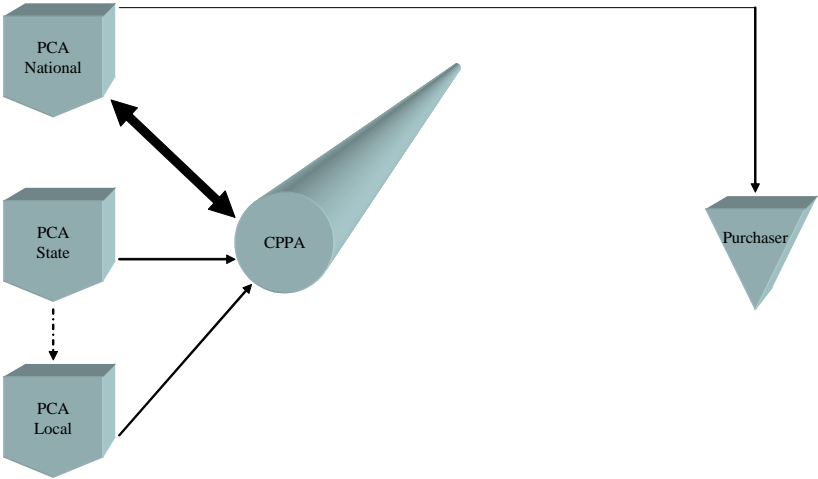
the relationship of the affiliates to each other.[9] Vertical public purchasing agreements are the legal covenants governing cooperative public purchasing where at least some of the members to the agreements (PCAs) are hierarchically ordered. This means that one PCA possesses the ability to limit another PCA from participation. An example is an agreement between a state and its localities such as that shown here. In this case, the state can compel local governments under its jurisdiction to purchase from the CPPA entity, and thus can dictate to the locality PCAs what it may and may not purchase. Under these circumstances, the muddle of common agency is amplified to further diminish transparency and accountability.

However, vertical agreements are different from horizontal public purchasing agreements. The latter CPPAs describe the legal covenants for member PCAs that are the same type of legal entities or when member PCAs have the same legal status. Examples include a CPPA with PCAs composed of states without their localities, or a CPPA with both states and the national government comprising the PCAs. This means that since each PCA has equal legal stature, the likelihood of coercive enjoinment is diminished. Thus, in vertical CPPAs, one or more members to the agreement have a subsidiary legal status inferior to other members and may be forced to join against their will. In Figure 2, this is depicted by the locality PCA being inferior to the state PCA.

So what does this demonstrate? Focusing on just the vertical relationship, the Buyer model implies that the dictatorial aspect of vertical public purchasing agreements may harbor the potential for distortion of purchases carried out by subservient PCAs. Consider that because cities and counties often wish to secure local suppliers to help in carrying out policies related to economic development, minority preferences, and economic efficiencies, a mandate to partake in a vertical public purchasing agreement might have the effect of distorting local purchasing contracts, especially toward larger suppliers usually found at a higher governmental level. Although there might be economic gains captured through economies of scale, the distortion of these other local goals may be unwanted and in fact might produce significant unintended consequences that remain invisible or negligible at the higher level. Thus, the Buyer model can help determine rivals, subsidiaries, and mutual partners. However, as Figure 3 will show, it is not the only theoretical paradigm of cooperative public purchasing.

**FIGURE 3**

## Piggyback Model of Cooperative Public Purchasing



## Piggyback Model

A second option for government entities to engage in CPP is to procure through an existing public cooperative affiliate (PCA) that has the means or expertise to buy in bulk and will then negotiate contracts for the other PCAs who are parties to the same agreement.  This relationship is depicted in Figure 3.  Prominent examples of this are when localities utilize interlocal agreements or when states use the U.S. federal GSA schedule.[10]  The GSA schedule is an internal buying organ of the federal government for its departments and agencies.  The GSA also allows states, localities, tribes and other entities to use its purchasing schedule.  In this way, a public cooperative affiliate is used by other parties to the cooperative purchasing agreement when they use a piggyback approach to obtain goods and services.

There are three important details in the Piggyback model depicted in Figure 3.  The first refers to the lack of a CPPA enterprise like that found in the Buyer model.  This means that the cooperative public purchasing agreement has no separate administrative organization outside of the member PCAs. A second point about the Piggyback model centers on the relationship between the national PCA and the other PCA members.  Notice that the linkage between the national PCA and the CPPA is bi-

directional. This suggests the possibility that the national PCA might simultaneously be both a principal for the CPPA (as a member) and an agent of the CPPA (as a buyer). In turn, this may imply that the incentives of the principals and agents in this model may be exhaustively aligned. Outside the scope of this chapter, such a possibility is a clarion call for future research in this area.

A third point to make about the Piggyback approach concerns the possibility that other PCAs could buy piggybacked depending upon the commodity or service being purchased. Consider the feasibility that although Figure 3 indicates that the national PCA is the only purchaser, there is the potential for conceiving of a meta-Piggyback model where each PCA brings unique commodity or service technocratic expertise to the group, which results in economies of scale such as those found in electrical cooperatives. Thus, the simultaneity depicted between the national PCA and the CPPA shown here might apply to other PCAs as well. This leads to the final model of cooperative public purchasing.

## Broker Model

Agency theory prompts consideration of the hypothetical case when a cooperative agreement exploits an organization external of the CPPA structure to make purchases for the membership. Like the Buyer and Piggyback models depicted in Figures 2 and 3, a Broker agreement (shown in Figure 4) can also have PCAs that are vertically or horizontally ranked as denoted by the dotted line from the state PCA to the local PCA. However, the layers of agency make clear that the organ directly charged with purchasing for the group is outside the direct control of the CPPA and its enterprise, and this may introduce principal-agent pathologies such as moral hazard and higher information asymmetries that may be lower in the Buyer or Piggyback paradigms.

Leaving the issue of pathologies aside for the moment, the Broker model appears to imply a potential for issues which might emanate from the contractual flow of goods and services, issues of ownership, and how the purchaser obtains title of the goods and services.

Dependent upon the role of the Broker in these activities, the consequences could be momentous. For example, if the Broker takes title first and acts as an agent of the PCAs, in a very real sense, the Broker acts as an indirect agent by selling to the PCAs.

**FIGURE 4**
**Broker Model of Cooperative Public Purchasing**



In contrast, under a framework of direct agency, the PCAs who are not the Broker contract directly with the supplier, but on the basis of prices, specifications, and terms negotiated by the Broker. This means that the PCAs contract with, and take title directly from, the supply source. Under these circumstances, the Broker merely acts as a buying agent for parties to the cooperative agreements and to cover its costs, may obtain a commission for its services which might cover the staff and operations budget of the Broker and the outlays associated with the cooperative public purchasing agreement (CPPA) and its enterprise.[11]

Although these considerations may also be present in the Buyer and Piggyback models, the Broker system highlights the options available and their potential consequences for adopting one method of purchasing over another. The Broker model also highlights other issues. For example, it is not self-evident why a government entity would join a CPPA that utilizes a Broker. Perhaps the answer may be that the CPPA can provide other benefits to its affiliates beyond lowest price (e.g., expediency, political neutrality, and networking). Indeed, even a private for-profit Broker may be a powerful organization, which through experience, can leverage the purchasing power of its customer CPPAs. Moreover, it is the private Broker that nurtures vendor relationships, provides expertise, streamlines the bidding and purchasing process, and probably provides a single contract to its customers. If these obligations

were to be carried out under Buyer or Piggyback conditions, the purpose of the CPPA (networking) may in fact be mitigated. For these and other reasons, CPPAs that utilize a Broker can be justified on defensible grounds, but only when there is a clear goal that is being met. That goal can be delineated by agency theory.

Indeed, there appears to be an example of the Broker model in Novation, LLC., which is a supply services company to the health care industry based in Irving, Texas. Novation has more than 2,400 "members" that include, among others, an alliance of not-for-profit hospitals and academic health centers. Another example of this model might be U.S. Communities. U.S. Communities bills itself as a "Purchasing Alliance" that is a nonprofit public benefit corporation whose sponsors include various associations and extra-governmental organizations like the National Institute of Governmental Purchasing, National Association of Counties, the National League of Cities, and the United States Conference of Mayors. As a broker, U.S. Communities acts as an agent for the particular state and local governments purchasing through the alliance, yet the state and local governments have no direct impact on how things are purchased through the alliance.

## DEFINITION OF COOPERATIVE PUBLIC PURCHASING

Having reviewed the three agency models depicted in Figures 2 through 4, one may be left wondering, what exactly is a public cooperative? Do the PCAs, the CPPA, the CPPA enterprise, the Broker, the purchaser, or all of these together comprise a public purchasing cooperative arena? Can one or more elements be excluded while still remaining true to the concept of a public purchasing cooperative arrangement? Conceptually, Figures 2, 3 and 4 offer some insight to answer these questions by specifying the chain of agency linkages from PCAs to the final purchaser for the CPPA. It is clear from the three models that in general, the common elements of all public purchasing cooperatives are the PCAs and the CPPA, but the departures between the models lead to the following comprehensive definition of a public purchasing cooperative:

> a public purchasing cooperative consists of a collaborative agreement between two or more governmental entities that funnel organizational and monetary resources into a purchasing syndicate which guides, regulates, and sanctions the conduct of the cooperative purchasing agent. Membership may be

voluntary or compulsory based upon the horizontal or vertical relationship of the affiliates.

To many this may be seen as merely re-describing *examples* of public cooperatives, especially in Figures 2 through 4. However, this is not the case for a very substantial reason.

Consider that knowledge is generated either through inductive or deductive processes. Because of this, some readers may validly believe that the term "cooperative public purchasing" should be clearly conceptualized early in any research endeavor, but that would be committing what has been claimed from the outset, namely, both *observation and theory* should drive the understanding of cooperative public purchasing. To rely only on observation is to remain normatively prescriptive without rational or logical justification. Thus, agency theory is the metaphorical glue that holds together the elements of cooperative purchasing by offering a framework for analysis, interpretation, and definitional clarification.

## DISCUSSION

Each of the models, especially Figure 1, treats some elements as "black boxes," (e.g., government entities are comprised of the purchasing departments and staffed by procurement officials). Using both inductive and deductive means, this chapter identifies crucial concepts and how they relate by presenting three models of cooperative public purchasing. Agency theory provides a framework to model the stakeholders in the cooperative public purchasing process, as well as introducing key terms that can be used by practitioners and academics alike. It is hoped that the models outlined here will lay the groundwork for future theoretical and practical work which might lead to a greater understanding of incentives for potential gains and hazards of engaging in specific types of cooperative public purchasing arrangements. As a result, cost-benefit analyses will be enhanced, leading to more effective long-range strategic planning by purchasers.

Developing a theory in cooperative public purchasing is important for several reasons. The theory elaborated here conceives of purchasing in a new and different way which can explain some counterintuitive incentives faced by public purchasing stakeholders. Because of this, theory can help to predict the behavior of stakeholders in terms that make rational assessment easier. In turn, this will clarify the needs and

goals of government entities in a way that encourages governments and other public organizations to design and adopt the most beneficial types of cooperative agreements given their economic and political needs in procurement. Obviously, this leads to a research agenda which seeks to untangle under what circumstances the models outlined here are useful.

One of the most surprising findings of this study was the realization that there is a chain of agency involved in any public purchase. The existence of these layers of agency might explain why public purchasers often feel pulled in different directions by trying to serve multiple masters. Unlike purchasers found in the private sector, public purchasers face a myriad of divided loyalties based upon the presence of both immediate and extended principals whose goals are often in conflict. Sensing that they have divided loyalties, they feel caught between competing demands for their time and efforts, yet they do not fully understand their predicament nor do they have a solution to this malaise. One should remember that in the private sector, there is one over-riding goal, and that is to maximize profit. However, in the public sphere, no such fundamental objective is clear in every case. To take but one example, "good government" can mean many different things to different principals, and to a great extent, the ambiguity of how to operationalize "good government" in public purchasing remains elusive. This chapter presents a theory that brings about a greater awareness of this situation so that purchasers might become more efficient in serving their stakeholders in and out of government.

Applying agency theory to CPP allows one to see the potential for adverse selection in several areas. Returning to Figure 1, due to informational asymmetries, there can be the wrong candidates elected to public office due to the rational ignorance of voters and the superficial campaign techniques so prevalent today (Prier, 2003). In turn, they may advocate the wrong procurement policies. It may be that procurement officials with the wrong sets of skills to knowledgeably engage in CPP populate purchasing divisions. If they have the right skill sets, they may be utilizing the wrong CPP for the objectives of their principals. All of these instances flow from an understanding of adverse selection embedded in agency theory.

Agency theory may also hold promise for modeling intergovernmental contracting, which is the leveraging of assets by cooperating with public agencies to provide goods and services to constituent end-users. Because a fee is paid by one government to

another for the provision of services, it might be considered a form of cooperative purchasing (possibly a Piggyback).  If this is the case, what are the advantages and traps of thinking about intergovernmental contracting as a cooperative agreement?  Furthermore, what are the legal and ethical implications for structuring a cooperative agreement in this way?

There are other questions evoked by this study.  What are the motivations for joining cooperative public purchasing agreements? Is it merely economies of scale, or is it the opportunities to network or streamline administrative functions? Is there an impact on individual member affiliates if the agreement allows members to choose goods and services a lá carte?  Furthermore, are specific cooperative agreements better for some entities than others, or is the choice of appropriate agreements contingent upon individual circumstances? Indeed, what characterizes a good cooperative, and under what circumstances are various cooperative models functionally appropriate?  In practice, are principal-agent pathologies endemic to one model and not others? If one believes that cooperative arrangements should follow many of the principles suggested by both the International Association of Cooperatives and the National Business Cooperative Association, there are potential violations of these principles in practice. They may be justified, because there may be reasons for exceptions.  For example, although cooperative enterprises are believed to require ownership by independent and autonomous members with membership being open and voluntary, vertical relationships among PCAs may purposefully violate these expectations depending on the goals of each PCA.

The schema presented here presses organizations to decipher their immediate and long-term goals.  Consider that at this juncture, it is unclear whether joining a cooperative serves an operational or strategic function or both, and without a substantive rule to judge the validity of any of these assertions, one cannot identify the differences between behaviors that are operational or strategic. This chapter is a first step in that direction.

It should be remembered that testable hypotheses are needed to determine if purchasers' decisions are aligned with policymakers' desires, and without a theory on which to base the propositions, practitioners and academics are left to random claims of descriptive tendencies.  Whether or not the theory and models offered here are valuable is left to the reader to decide.  However, they are proffered with

the strong conviction that to continue the trend toward cooperative purchasing without theoretical direction is tantamount to making purchasing decisions based on blind faith.

## ACKNOWLEGMENTS

## NOTES

1. The discussion here relies heavily on the review articles by Kiser (1999) and Petersen (1993).

2. Although the literature has devoted substantial effort in understanding transaction cost economics (see Williamson [1985]), the focus in the present analysis remains concentrated on the principal-agent perspective (e.g., Milgrom and Roberts 1992, esp. Chs. 5-6).

3. Moral hazard refers to the principal's increased risk of suffering negative consequences resulting from problematical behavior of the agent. It is present because the agent may benefit from the outcome or will not suffer the adverse consequences of her own behavior.

4. Although the principal may be able to observe the agent's action in some circumstances, the observation typically requires costly monitoring. Monitoring might obtain information on the agent's ability, carefulness, laziness, reliability, and trustworthiness, to name a few characteristics.

5. For some organizational pathologies in the purchasing supply chain, see Mishra, Heide, and Cort (1998).

6. Note that purchasing associations themselves can form a separate organization for cooperative purchasing.  For example, educational cooperatives and political units such as county education offices have formed and are now member PCAs in the Association of Educational Purchasing Agencies. It is organized through a Memorandum of Understanding between all participating states.

7. The node defined as the purchasing department could also be the agency allowed to purchase in a decentralized structure, but for the purpose of exposition, the analysis is simplified.

8. Although there are numerous examples of this and other CPPAs around the world, the current analysis uses the U.S. to simplify the discussion. For other examples in the U.K., see Gershon (2004).

9. For clarification purposes, only the vertical possibility is portrayed in all figures presented here.

10. The GSA is actually a vertical agreement because all of the parties (PCAs) to the agreement are not between the same caste of legal entities.  In other words, because the GSA allows both localities and their states, and because the state PCA has the unilateral ability to create and abolish the locality PCA, this is a vertical CPPA.

11. This would be analogous to outsourcing the purchasing function.

12. See Prier (2003) for an extended critique of divided loyalties.

13. For example, a state may require a locality to join and use a particular CPPA in order to get best price.  In this case, other local considerations are trumped by the requirement.

## REFERENCES

Anderson, R, & Macie, K. E. (1996). "Understanding the Differences in Group Purchasing." *NAPM InfoEdge*, 2 (4). [On-line]. Available at www.ism.ws/pubs/InfoEdge/InfoEdgearticle.cfm?ItemNumber=101 59 (Accessed on July 1, 2004).

Aylesworth, M. M.  (2003). "Purchasing Consortia in the Public Sector Models and Methods for Success."  Paper presented at the 88th Annual International Supply Management Conference, May 18-23, Nashville, Tennessee. [On-line]. Available at www.ism.ws/ResourceArticles/ Proceedings/2003/AylesworthJI.pdf. (Accessed on July 2, 2004).

Becker, G. S. (1983). "A Theory of Competition among Pressure Groups For Political Influence." *Quarterly Journal of Economics, 98*: 371-400.

Coase, R. (1937). "The Nature of the Firm." *Economica, 4*: 386-405.

Demsetz, H. (1983). "The Structure of Ownership and the Theory of the Firm." *Journal of Law and Economics, 26* (2): 375-390.

Demsetz, H. (1971). "On the Regulation of Industry: A Reply." *Journal of Political Economy, 79*: 356-363.

Dewatripont, M., Jewitt, I., & Tirole, J.. (1999). "The Economics of Career Concerns, Part II: Application to Missions and Accountability of Government Agencies." *Review of Economic Studies* (Special Issue: Contracts) *66*: 199-217.

Dixit, A. (2002). "Incentives and Organizations in the Public Sector: An Interpretative Review." *Journal of Human Resources, 37* (4): 696-727.

Faith, R. L., Leavens, D. L., & Tollison, R. D. (1982). "Antitrust Pork Barrel." *Journal of Law and Economics, 25*: 329-342.

General Accounting Office (1997, February). "Cooperative Purchasing: Effects are Likely to Vary Among Governments and Businesses" (GAO/GGD-97-33). [On-line]. Available at www.gao.gov/archive/1997/gg97033.pdf. (Accessed on July 4, 2004).

Gershon, P. (2004). *Releasing Resources to the Front Line: Independent Review of Public Sector Efficiency*. Crown, UK: Add Name of publisher.

Goodwyn, L. (1976). *Democratic Promise: The Populist Movement in America*. New York: Oxford University Press.

Grossman, G. M., & Helpman, E. (1996). "Electoral Competition and Special Interest Politics." *Review of Economic Studies, 63*: 265-286.

Jensen, M. C., & Meckling, W. H. (1976). "Theory of the Firm: Managerial Behavior, Agency Costs, and Ownership Structure." *Journal of Financial Economics, 3*: 305-360.

Johnson, R. H. (1983). "The New Populism and the Old: Demands for a New International Economic Order and American Agrarian Protest." *International Organization, 37*: 41-72.

Kiser, E. (1999). "Comparing Varieties of Agency Theory in Economics, Political Science, and Sociology: An Illustration from State Policy Implementation." *Sociological Theory, 17* (2): 146-170.

Knapp, J. G. (1969). *The Rise of American Cooperative Enterprise (1620-1920)*. Danville, IL: The Interstate Printers and Publishers.

Lloyd, R. E. (2000). "Government Contracting Pathologies." *Acquisition Review Quarterly* (Summer): 245-258.

Lohmann, S. (1998). "An Information Rationale for the Power of Special Interests." *American Political Science Review, 92*: 809-827.

Milgrom, P., & Roberts, J. (1992). *Economics, Organization, and Management*. Upper Saddle River, NJ: Prentice Hall.

Miller, M. G. (1937). "The Democratic Theory of Cooperation." *Annals of the American Academy of Political and Social Science, 191*: 29-37.

Mishra, D. P., Heide, J. B., & Cort, S. G. (1998). "Information Asymmetry And Levels Of Agency Relationships." *Journal of Marketing Research, 35*: 277-295.

O'Toole, L. J. Jr. (1995). "Diversity or Cacophony? The Research Enterprise In Public Administration." *Public Administration Review, 55* (3): 293-297.

Petersen, T. (1993). "The Economics Of Organization: The Principal-Agent Relationship." *Acta Sociologica, 36*:277-293.

Prier, E. (2003). *The Myth of Representation and the Florida Legislature: A House of Competing Loyalties, 1927-2000*. Gainesville, FL: University Press of Florida.

Ramsay, J., & Caldwell, N. D. (2004). "If All You Have Is a Hammer, Everything Looks Like a Nail: The Risks of Casual Trope Usage in Purchasing Discourse." *Journal of Purchasing and Supply Management, 10*: 79-87.

Ross, S. A. (1973). "The Economic Theory of Agency: The Principal's Problem." *American Economic Review, 63* (2): 134-139.

Saloutos, T., & Hicks, J. D. (1951). *Agricultural Discontent in the Middle West 1900-1939*. Madison, Wisconsin: University of Wisconsin Press.

Taylor, C. C. (1953). *The Farmers Movement 1820-1920*. New York: America Book Co.

Telser, L. G. (1969). "On the Regulation of Industry: A Note." *Journal of Political Economy, 77*: 937-52.

Terry, L. D. (1999). "From Greek Mythology To The Real World Of The New Public Management And Democratic Governance (Terry Responds)." *Public Administration Review, 59* (3): 272-277.

Tulip, S. (1999, December). "Joining Forces." *Supply Management*: 24-30.

Williamson, O. E. (1985). *The Economic Institutions of Capitalism*. New York: The Free Press.

Wooten, B. (2003). "Cooperative Purchasing in the 21st Century." *Inside Supply Management* 14: 4-6.

# Organisation Theory: The Principal-Agent Perspective

By Jan-Erik Lane

*Abstract-* Today much relevant questions concern Who get what, when and how?, due to the incredible rise in the remuneration of the economic, cultural and political elites in the large organisations around the world. A suitable conceptual framework for the analysis of the fundamental question, namely Cui Bono?, is the principal-agent approach from recent advances in game theory. The skyrocketing of the salaries and bonuses of CEO:s in the private sector and the spreading out of corrupt practices in the public sector forces the social science to ask the quid pro quo question about the relationship between the remuneration of agents and their delivery of outputs to the principal. It is truly fruitful for the understanding of political organisation in whatever form it takes. Politics everywhere is about contracting, introducing a web of contracts between principal and agents. The shape of these contracts determines the real constitution of a country.

*Keywords:* organisation theory, incentives, contracting, considerations in contracts, quid pro quo, cui bono, asymmetric information, simple contracts – complex organisation, political organisation: demos versus politicians and officials.

ORGANISATIONTHEORYTHEPRINCIPALAGENTPERSPECTIVE

*Strictly as per the compliance and regulations of:*

# Organisation Theory: The Principal-Agent Perspective

Jan-Erik Lane

*Abstract-* Today much relevant questions concern *Who get what, when and how?,* due to the incredible rise in the remuneration of the economic, cultural and political elites in the large organisations around the world. A suitable conceptual framework for the analysis of the fundamental question, namely *Cui Bono?,* is the principal-agent approach from recent advances in game theory. The skyrocketing of the salaries and bonuses of CEO:s in the private sector and the spreading out of corrupt practices in the public sector forces the social science to ask the quid pro quo question about the relationship between the remuneration of agents and their delivery of outputs to the principal. It is truly fruitful for the understanding of political organisation in whatever form it takes. Politics everywhere is about contracting, introducing a web of contracts between principal and agents. The shape of these contracts determines the *real constitution* of a country.

*Keywords: organisation theory, incentives, contracting, considerations in contracts, quid pro quo, cui bono, asymmetric information, simple contracts – complex organisation, political organisation: demos versus politicians and officials.*

## I. Introduction

The principal-agent model offers yet another framework for analysing the organisation of human activities(Ross, 1973; Grossman and Hart, 1983;Sappington, 1991;White, 1992; Ackere, 1993; Althaus, 1997). Its strength is that it underlines incentives more than rules as in many organisation approaches. The focus is upon the *web of contracts* that link people together in an organisation, analysing them with the newly developed concepts of in the economic theory of information (Bircher and Butler, 2007).

Looking at relations between actors as contractual links between principals and agents has proved insightful with regard to understanding employment/sharecropping in agriculture, the work of attorneys, the doctor-patient relationship or investor-broker interaction as well as the entire business of insurance. Yet, there has been great reluctance to apply the principal-agent model to politics, because the key concepts do not seem to capture the essence of politics in well-ordered societies, namely to safeguard the national interest or common good of citizens.

Human beings have developed great skills in organizing activities so that an ever increasing output of goods and services is possible. Thus, organisations of various kinds play a major role in social life every day. Organisation theory and management studies have contributed lots of studies with numerous insights into the operations of organisations, market based as well as non-market organisations. This intense research has resulted in a number of theoretical approaches. These frameworks underline a variety of factors in or aspects of organisations: e.g. planning, strategy, internal organisation – external relations, hierarchy, division of labour, bounded rationality and institutionalisation.

The aim of this paper is to raise the question *CUI BONO?* in relation to organisations. It is hardly an exaggeration to say that organization theory and management approaches have been much concerned with efficiency, meaning the successfulness of the organization. Also the big branch of organization studies that deny the possibility of efficiency is occupied with the same perspective: outputs, outcomes, resources, strategy, leadership, etc., although underlining the relevance of so-called garbage can patterns of organization and management. The quest for efficiency of organizations in both classical management theories and public administration approaches and its rejection in the bounded rationality perspective upon organizations, launched by H. Simon and J. March, has resulted in an intense debate about the nature of organizations and the limits of management. But neither of these two theoretical perspectives entails much for the crucial question about organisations, namely: *Cui bono*? Even the most radical approach to organisation, denying completely the relevance of concepts like effectiveness and productivity to understand real life management, preferring to talk about organised foolishness, myths and institutional legacies (Olsen, 2010: Brunsson, 1985), does not touch the fundamental *Quid pro quo* questions in organisations: Who gains?

Interestingly, the rational choice approach in the social sciences has been accused of being linked logically with the efficiency focus. If people are summed to act so as to maximise their goals in a rational manner, then arguable they would do the same when managing organisations. However, the entailment does not hold. The management of an organisation involves collective decision-making among a group of people – the managers. Each of them may pursue their goals according to the requirements of individual rationality, yet when combined these individual decisions may lead

*Author: An independent scholar, professor at three universities.*
*e-mail: janeklane@googlemail.com*

to suboptimal decision-making and even chaos or foolishness.

The quid pro quo question in relation to organisations leads to the emphasis upon contracting, asking the following: What have people agreed upon to? Against what pay? With what effort? How are the outputs to be measured? And what is involved in the evaluation of performance: firing, bonus, new contract, etc? The content of any contract is its consideration, meaning the expectations that the parties bring to the agreement. The organisation is a WEB of contracting and management is the handling of these contracts, from the beginning – ex ante – to its fulfilment – ex post.

Studying organisations as webs of contracts and their management, the principal-agent framework from recent advances in game theory appears most promising. Thus, we ask:

1. Why it is easy to organise lots of taxi services in a huge capital like Yangon?
2. How come the remuneration of CEO:s is out of hand?
3. How can politicians become superrich?

### a) The Stylised Principal-Agent Model

According to Rasmusen (2006), the principal-agent model includes a principal searching to maximise the value of some output(s) V by means of contracting with a set of agents, remunerating them for their efforts in producing the output. The payment of the agents derives from the value of the output of the agents, meaning that the principal-agent contract must involve considerations covering the *ex ante* to the *ex post* stages. With a considerable time lap between the making of the contract and the fulfilment and its evaluation, problems of asymmetric information and transaction costs arise (Rao, 2002).

The principal-agent framework has enjoyed far reaching success in modelling interaction between persons where one works for the other. This interaction is to be found in many settings, such as agriculture, health care, insurance and client-lawyer (Ross, 1973; Rees, 1985: Laffont and Martimort, 2002). As a matter of fact, the principal-agent problematic is inherent in any employment relationship where one person works for another, who pays this person by means of the value of the output. Whenever people contract with others about getting something done, there arise the typical principal-agent questions:

1. What is the *quid pro quo* between the principal and the agent – the contractual considerations?
2. How can the principal check the agent with regard to their agreement – the monitoring problem?
3. Who benefits the most from the interaction between principal and agent – who takes the surplus?

These questions concerning principal-agent interacting arise whenever there is a long-term interaction between two groups of people, involving the delivery of an output against remuneration as well as a time span between the making of the contract and the ending of the relationship with the final delivery of the output, Let us apply this conceptual framework to three kinds of organisation in order to demonstrate that it illuminates the pattern of interaction.

## II. Taxi Services in Yangon: The Principal on Top

Powerful forms of connecting people may result from very simple contracts between principals and agents, like in sharecropping. They may last long and need not even be formalised in written agreements. They may involve hundreds of people working as agents for one single principal, owning the assets involved in the production of services.

### a) Taxi Organisation:

1. *Principal:* Owner of the cars, with goal to maximise profits from taxi services;
2. *Agents:* Renting the car for 12 $ a day with a guarantee of 300 $ for damages as first down payment. All running costs are born by the agent and the car is checked in detail at every round of contracting period.

*Outcome:* The principal, who is risk avert, provides the car but the agent has to pay all repairs, either with the down payment or additionally through a loan from the principal. The agent will drive the care as long as he/she can raise every day > 12 $ plus the running costs and the repair costs. This contract is attractive for people whose reservation salary is very low or zero. It is also incentive compatible, as the driver gains more by being active. This organisation tends to be stable. Since unemployment is high in Yangon, the remuneration of agents can be kept as low as possible, securing a nice profit to the principal, who bears little risk.

## III. The Joint-Stock Organisation: Agents on Top

Besides the trillions of daily on-spot contracts in the markets, there occur several forms of principal-agent contracting, introducing organisation into social life. A simple principal-agent contracting was described above, but there are others forms than one to one, like one to many, many to one and many to many. In the private sector, firm organisation varies from small partnerships to giant enterprises with more than one hundred thousand employees. It is all based upon contracting between principals and agents, which is why law and lawyers loom so large, i.e. private law.

### a) Firm Organisation:

1. *Principal:* Owners of the shares: a few big owners plus an ocean of small owners with the goal of maximising the value of their holdings of stock;

806

2. *Agents:* The CEO:s, who are risk avert, receiving a fixed salary plus yearly bonus, decided usually at discretion. The CEO can be fired at any moment but receives a so-called golden handshake. He/she employs the other employees on standard wage contracts – internal organisation – or on the basis of outsourcing. All the agents are paid by means of the market sales of the output of the firm, where the CEO:s maximise their remuneration in total.

*Outcome:* The owners will need lots of monitoring to find out what is going on and whether the CEO:s make an effort. Thus, they wish to list the firm on the bourse, harbouring instantaneous evaluation. The risk of the owners is the occurrence of asymmetric information, both ex ante (adverse selection)and ex post (moral hazard). This organisation tends to be unstable, as the CEO:s manage to use various strategies to push up their remuneration almost to the level of looting.

The instability in the firm organisation shows up in the constantly increasing remuneration packages of the CEO:s, where the spread to other employees have multiplied several times during the last 50 years. This is true of both the fixed salary and the yearly bonuses, which tend to be paid more or less automatically. It has happened that bonuses become permanent remuneration whatever the result of the firm is.

Neither economic decision theory nor management theory has any clear explanation of the tendency of the CEO:s to prevail to significantly in the firm organisation. The only credible explanation is that shareholders are easily manipulated by the CEO:s due to the enormous asymmetric information plus the large room for the CEO's to enter collusion by making coalitions with board members, like first and foremost the chairman of the board of the company. As effort is not observable and costly to enforce, shareowners chose to believe in the story of the CEO, often until it is too late.

There is no remedy to this advantage of the agent. Making the CEO part owner of the firm has been proposed but the future price of his stock options tends to be set extremely low. A radical solution is that the big owners become the CEO:s, but this is only feasible for some firms, like e.g. HM.

The remuneration of CEO:s could skyrocket when various forms of commissions are added to the salary, for instance when company activities are sold or bought. The remuneration of the CEO of NOKIA before it was sold to MICROSOFT is an excellent example. Firms that are owned by consumers themselves, like COOP, are exceptionally vulnerable to the claims of CEO:s, when excessive.

It is an often debated fact that the total remunerations of agents has gone up astronomically over the last decades in the firm organisation, resulting in rapidly increasing inequality in both Western societies and Eastern or South East Asia. The basic reason is hardly a shortage of CEO:s or a dramatic increase in management skills, but simply the instability inherent in the principal-agent interaction in firm organisation due to asymmetric information. When the CEO:s are hired, there is the adverse selection problem of failing to recognize pretending and when they have been hired, there is the moral hazard problem of shirking. The shareholders are so afraid of these two major difficulties in firm management that they are prepared to throw almost any amount on money upon them. It has happened that the CEO:s capture almost all the profits of a joint-stock company in the form of bonuses: It would be better for its shareholders to sell this company (*Husqvarna*) to these CEO:s! Public joint-stock companies with the state as the owner are run with the same principal-agent interaction: the CEO agents on top. The process of incorporation all of Europe has resulted in huge increases in their remuneration, like Swedish *Vattenfall*.

## IV. The Remuneration of Politicians

Political science teaching often starts with the observation that roughly 50 per cent of all existing countries today have a democratic regime of some sort while the rest of the countries either are authoritarian regimes or so-called failed states, i.e. countries in anarchy. This distinction between democracy and non-democracy has been a very central research topic since after the Second World War, especially as the number of democracies has increased during the last decades. A large number of factors have been examines, exogenous as well as endogenous ones, like the economy, social structure, ethnicity, religion, openness, historical legacies, etc.

A completely different way of approaching this research issue in the social sciences, economics and politics is to start from the *quid pro quo* question. In non-democracies, the remuneration of politicians tends to be much higher than in democracies. And in failed states, the predicament of anarchy opens up for the looting strategy, which may pay off handsomely for rebels, jihadists and drug traffickers. In kingdoms or sultanates, the existence of patrimonial authority implies that imperium and dominium, public authority and private ownership are fused. Thus, e.g. the Saudi family is the owner of the oil riches of the country. Moreover, the sultan of Oman *Qaboos bin Said Al Said* receives all state revenues as his, thereafter writing checks to the public budget, as signs of generosity.

In authoritarian one-party states, the political leadership forms a most wealthy click, like in the Khanates and China. Why start a transition to democracy when so much of wealth is at stake for the economic fortunes of the rulers? In his detailed enquiry in the fate of African states after the coming of

47

independence from the Europeans, British historian Meredith documents an almost incredible list of rulers who enriched themselves through embezzlement. No wonder that many of them attempted to stay on as long as possible, even for 2-3 decades! The political agents will try to capture as much as possible of the value V of the game, i.e. the country GDP, unless hindered by competing agents or guardians like courts or the Ombudsman (Public Protector).

*b) Political organisation:*

*Principal: demos*, citizenry, electorate, population

*Agents:* politicians, parties, legislators, judges, Ombudsman, bureaucrats, officials, agencies, boards, etc.

*Incentives:* What drives the agents? And do they really improve for the principal?

The state is a much more complicated organisation than the firm. It likewise involves lots of laws and regulations, i.e. public law. Perhaps this is why the principal-agent approach has not been applied systematically? In any case, one needs to ponder on how the interaction is to be modelled with the variety of players. Principal-agent interaction in constitutional democracies is very different from that of non-democracies. A number of models have been launched: Barro, 1973; Ferejohn, 1986; Weingast, 1989; Rao, 2002; Besley 2006; Helland and Sörensen, 2009.Yet, the central question is the following: How do constitutional democracies reduce the upper hand situation of political agents in non-democracies?

## V. Constitutional Political Organising

The following assumptions appear the most likely to be adequate for modelling principal-agent games in a constitutional democracy:

1. The principal of the democratic state is the demos, or the electorate – *body politic*;
2. The set of political agents covers three groups: governments and its bureaucracy, the legislators and the judiciary – *trias politica*;
3. Politicians offer the voters alternative policy packages about how the state may improve upon society, or total value V;
4. The remuneration of the political agents are separated from the resources of the *fiscus*, the state coffers;
5. The remuneration of politicians is fixed, including pensions, in order to avoid the appropriation of the *fiscus*;

These maxims of constitutional democracy seem enough to introduce the distinction between the public and the private, which was so confused in all forms of oriental despotism, as well as solve the *appropriation problem* in politics and public administration, as Max Weber conceived it (Weber,

1978). The modern bureaucracy and its superior performance to patrimonial administration is only feasible when officials are paid predictably, meaning that they are little incentive to *appropriate* the recourses of Bureaux or engage in looting in society.

However, we need a few more maxims:

6. The principal will only be able to control the set of political agents when they are set in competition with each other;
7. Political competition is as vital to democratic politics as firm competition is to the market;
8. Political competition favour the interests of the demos, pitting the three branches of constitutional government against each other;
9. Political entry in competition must be open so that the authoritarian politicians cannot exercise political monopoly;
10. The judiciary operates on the principles of due process of law, to be found in either Common Law or Civil Law.

In order to tame the political agents and diminish their advantages, the principal has supported the evolution of distinct institutional mechanisms that restrain the political agents: viz. rule of law and the political market place. The hope is that the actions and decisions of politicians will enhance societal value, like for instance affluence and wealth.

## VI. Remuneration and Value in Principal-Action Games

It is an axiom in the principal-agent model that the agents are paid from the value of the output they deliver for the principal, who is the residual claimant. The principal wants to maximise that value, but he/she must present the agents with an incentive compatible contract, paying more for higher effort. As there is no guarantee that higher effort will actually be forthcoming or succeed in baking a bigger cake, principal-agent contracting is replete with failure, which could leave the principal pay all the value to the agent – the case of looting. In the worst case scenario, the principal pays for high effort but the agents employs the strategies of pretending and shirking to deliver a meagre output, resulting in a loss to the principal, as the value of the output does not cover the remuneration of the agents.

This is, of course, the fundamental *quid pro quo* problematic in all forms of contracting, private or public. In the organisation of taxi services above, the contract favour the principal, pushing the risk upon the agents. In firm organisation, it is the other way around. What about politics?

The state and political leadership concern an entire country, or nation, Thus, the value of the output of the political agents is their contribution to the total value in society, or the GDP. Moreover, the political agents are paid through taxes and charges upon the GDP. What is

the logic of the *quid pro quo* requirement, the consideration of the public contracts?

The most profound answer to this question is to be found in the theory of public finance, focussing upon the allocation to society of so-called public or semi-public goods (Musgrave and Musgrave, 1980). A country has a strong need for goods and services that are non-rival or non-excludable as well as joint in supply. As the market cannot supply these, only the public sector or the state can be relied upon. Market failure is the reason of the state.

Public or semi-public goods include law and order, peace and war, infrastructure, common pools, etc. In order to provide these services, political communities – governments at various levels – contract with a set of political agent to deliver them. What will be their remuneration for their achievements?

1. *Patrimonialism:* From the point of view of human known history, patrimonialism is the most frequently occurring structure of political leadership, at least until 1900. The remuneration of the political tends to go very high, at the same as there is constant struggle among contenders to the patrimonial assets. To stabilise the rulership, political leaders engages in huge aggrandizement project, which both deliver public goods and underlines their own position. When the subjugation of the principal, the population, becomes too excessive, spontaneous uproars follow, It takes a long for patrimonialism to accept the distinction between *crown* and *realm* – the so-called "*King's Two bodies*" (Kantoriwicz, 1957).At the core of all forms of patrimonialism whether in Europe, Americas, Africa or Asia is the consideration: How can the principal call upon the agents to deliver goods and services that further their interests, when opposition is met with arbitrary arrest, detentions and incarcerations?

2. *The authoritarian one-party state:* Patrimonilism (*l'etat c'est moi:* Louis XiV) was replaced by populist regimes that promised to fulfil the General Will of the principal, but in reality created the origins of totalitarian democracy (Talmon, 1952). Populist authoritarianism has occurred in several versions since 1800, but its key foundation is the manipulation of asymmetric information. The principal is deceived into supporting the agenda of the authoritarian elite by means of ideology and its myths: *France – la gloire* (Napoleon), the international proletariat (Lenin, Stalin, Mao), German nationalism (Hitler), Roman *grandeur* (Mussolini), Great Serbia (Milosevic), Kim Dynasty (North Korea), Zaire (Mobutu), etc. The remuneration of the authoritarian elite tends to be extremely high, including the taking of the babies of opponents (Argentina). Yet, the indirect costs may be much larger, as the rulers do not hesitate to put the entire country at risk. They may also be so cruel as to destroy society when threatened in power, like Mengistu in Ethiopia or Pol Pot in Cambodia. The authoritarian set of agents cannot accept any challenge from outsiders and does not hesitate to employ torture, sudden disappearances and assassinations to remove challengers or critiques.

3. *The Constitutional democracy:* To keep remuneration of political agents within reasonable bounds, the *quid pro quo* problem is here solved by very strict rules about the public budget - transparency. And to hinder that political elites replace their commitment to the welfare of the country with their own goals, there is detailed specification of rules of election and re-election – political markets. However, the direct and indirect costs of the politicians have certainly gone up in the last decade. Moreover, the costs of party operations keep escalating, creating a big grey zone where corruption may be suspected.

The indirect costs of the mistakes of political agents may be large, also in constitutional government. Thus, for instance the Bush family has born none of the misery that Operation *Cobra* (*Iraqi Freedom*) resulted in for ordinary people and servicemen. Now the Middle East is in total chaos: *bellum omnium contra omnes*. Similarly, the Putin policy against the Ukraine has proven very costly for the principal, the Russian peoples.

In the political markets, the costs of election may be extremely high in some countries. This is the problem of campaign fundsand its *quid pro quo*. Two questions: Can they be used as remuneration for the politicians? Do they involve a tacit contract to the effect that the politician (political party) is supposed to deliver outputs that favour the contributors (Peltzman, 1998)? The financing of the campaign expenses of political parties and individual politicians constitutes a grey zone between legality and corruption.

## VII. Conclusion

The principal-agent approach, developed in the economics of information and the game theory of successive moves in contracting (Rasmusen, 2006) may be employed to create a parsimonious theory of political organisation. It covers the essential aspects of principals versus agents, agent remuneration against the value of output to the principal, the monitoring of performance and conduct of political accountability as well as asymmetric information and its consequences for deception and manipulation.

## Literature

1. Ackere, A. (1993) "The principal/agent paradigm: Its relevance to various functional fields", European Journal of Operational Research, Vol. 70: 83-103.

49

2. Althaus, C. (1997) "The application of agency theory to public sector management", In G. Davis, B. Sullivan & A. Yeatman (eds,) The New Contractualism?, eds. Centre for Australian Public Sector Management, pp. 137–153.

3. Arrow, K. (1985)"The economics of agency" Pp. 37-51 in J. Pratt and R. Zeckhauser(eds), Principals and agents: The Structure of business. Boston: Harvard University Press,

4. Barro, R.J. (1973)"The Control of Politicians: An Economic Model," Public Choice 14 (Spring 1973): 19-42.

5. Besley, T. (2006) Principled Agents? The Political Economy of Good Government. Oxford: Oxford U.P.

6. Birchler, U. and M. Bütler (2007) Information Economics. London: Routledge.

7. Brazier, R. (1990) Constitutional and Administrative Law. London: Penguin Books.

8. Brunsson, N. (1985) The irrational organization: irrationality as a basis for organizational action and change. Chichester: Wiley.

9. Ferejohn, J.(1986)"Incumbent performance and electoral control", Public Choice 30: 5-25.

10. Ferejohn, J. and C. Shipa (1990) "Congressional Influence on Bureaucracy." Journal of Law, Economics, and Organization 6:1–20.

11. Furubotn, E.G. and R. Richter (2005) Institutions and Economic Theory: The Contribution of the New Institutional Economics. Ann Arbor: The University of Michigan Press.

12. Grossman, S. J., and O. D. Hart (1983) "An analysis of the principal-agent problem", Econometrica, Vol. 51: 7-46.

13. Helland, L. and Sørensen, R. J. (2009) "Hvorfor overlever politisk korrupsjon i representative demokratier?", Norsk Statsvitenskapelig Tidsskrift 2009, Vol 25. 3: 219-236.

14. Jowell, J. (1994) "The Rule of Law Today", in Jowell and Oliver, op.cit.:57-78.

15. Jowell, J. and Oliver, D. (eds.) (1994) The Changing Constitution. Oxford: Clarendon Press.

16. Kant, I. (1974) The Philosphy of Law. Clifton: Augustus M. Kelley.

17. Kantorovicz, E. (1957) The King's Two Bodies: A Study in Mediaeval Political Theology. Princeton: Princeton University Press.

18. Kelsen, H. (1961) General Theory of Law and State. New York: Russell & Russell.

19. Kelsen, H. (1967) Pure Theory of Law. Berkeley: University of California Press.

20. Laffont, J.J. and D. Martimort (2002). The theory of incentives: the principal-agent model. Princeton, New Jersey: Princeton University Press.

21. Lloyd, D. (1991) The Idea of Law. London: Penguin Books.

22. Lloyd, H. A. (1991) "Constitutionalism" in Burns and Goldie, (op.cit.), pp. 254-297.

23. Meredith, M. (1997) The State of Africa. London: Free Press.

24. Meredith, M. (2006) The Fate of Africa. Publicaffairs.

25. McIlwain, C.H. (1958) Constitutionalism, Ancient and Modern. New York: Cornell University Press.

26. Musgrave, R.A. and Musgrave, P. B. (1980) Public Finance in Theory and Practice. New York: McGraw Hill.

27. Neumann, F.L. (1986) The Rule of Law: political theory and the legal system in modern society. Leanington Spa: Berg.

28. Peltzman, S. (1998) Political Participation and Government Regulation. Chicago: University of Chicago Press.

29. Pennock, J. R. and Chapman, J.W. (eds.) (1979) Constitutionalism. New York: New York University Press.

30. Olsen, J.P. (2010) Governing through Institution Building. Institutional Theory and Recent European Experiments in Democratic Organization. Oxford: Oxford U.P.

31. Rao, P.K. (2002) The Economics of Transaction Costs. Basingstoke: Palgrave/Macmillan

32. Rasmusen, E. (2006) Games and Information: An Introduction to Game Theory. Oxford: Blackwell.

33. Reiss, H. (ed.) (1970) Kant's Political Writings. Cambridge: Cambridge University Press.

34. Rees, R. (1985) "The Theory of Principal and Agent", Bulletin of Economic Research, Vol. 37,1: 3-26.

35. Riley, P. (1983) Kant's Political Philosophy. Rowman & Allanheld Publishers.

36. Ross, Steven, (1973) "The economic theory of agency: The principal's problem", American Economic Review, 63(2): 134-139.

37. Sappington, D. (1991) "Incentives in principal agent relationships", Journal of Economic Perspectives 3(2): 45-66.

38. Schwöbel, C.E.J. (2011) Global Constitutionalism in International Legal Perspective. Leiden: Martinus Nijhoff.

39. Talmon, J.L. (1952) The Origins of Totalitarian Democracy. London: Secker & Warburg.

40. Tierney, B. (1982) Religion, Law, and the Growth of Constitutional Thought 1150-1650. Cambridge: Cambridge University Press.

41. Weber, M. (1970) Economy and Society. Berkeley: University of California Press.

42. Weingast, B. (1989) "The Political Institutions of Representative Government: Legislatures", in Journal of Institutional and Theoretical Economics, Vol 145: 693-703. Reprinted in Furubotn, E. and R. Richter (eds), The New Institutional Economics. (Tübingen: J.C.B. Mokr (Paul Siebeck) and college Station: Texas A&M Press, 1991).

50

43. Vile, M.J.C. (1967) Constitutionalism and the Separation of Powers. Oxford: Oxford University Press.
44. White, William D. (1992) "Information and the control of agents", Journal of Economic Behavior and Organization, Vol. 18: 111-117.
45. Wormuth, F.D. (1949) The Origins of Modern Constitutionalism. New York: Harper.

51

# CONDITIONAL ANALYSIS AND A PRINCIPAL-AGENT PROBLEM

JULIO BACKHOFF [1] AND ULRICH HORST [2]

ABSTRACT. We analyze conditional optimization problems arising in discrete time Principal-Agent problems of delegated portfolio optimization with linear contracts. Applying tools from *Conditional Analysis* we show that some results known in the literature for very specific instances of the problem carry over to translation invariant and time-consistent utility functions in very general probabilistic settings. However, we find that optimal contracts must in general make use of derivatives for compensation.

## 1. INTRODUCTION

In this article we analyze conditional optimization problems arising in the dynamic *Principal-Agent (PA) Problem* of delegated portfolio management. In these models, which belong to the class of contracting problems under moral hazard, an investor (the Principal) outsources her portfolio selection to a manager (the Agent) whose investment decisions the investor cannot or does not want to monitor.

Moral hazard problems have been first studied in [17, 18] in static environments and in [24, 23] in dynamic ones. In recent years, such problems have received renewed attention in the economics and financial mathematics literature. The seminal contribution [22] analyzed dynamic moral hazard problems in continuous time in which the output is a diffusion process with drift determined by the Agent's effort. The optimal contract, based on the Agent's continuation value as a state variable, was computed using sophisticated stochastic control and PDE methods. Using similar tools, [3] studied a PA model in which a risk-neutral Agent with limited liability must exert unobservable effort to reduce the likelihood of large but infrequent losses. In [25] a Stochastic Maximum Principle was applied to dynamic PA models to find first order conditions for optimality. In the most general case the Stochastic Maximum Principle leads to the characterization of optimal contracts through a system of fully coupled Forward-Backward Stochastic Differential Equations for which no general existence theory exists. These equations can typically only be solved explicitly when the analysis is confined to models driven by Brownian motion, specific preferences - typically linear, expected exponential or power utility functions - and information is symmetric, i.e. both parties observe the driving Brownian motion. We refer to the monograph [11] for a systematic survey of the mathematical literature on dynamic PA models and to [10] for a recent model of portfolio delegation under incomplete information which leads to even more complex dynamics.

Our work is motivated by that of Ou-Yang in [20] - which was in a sense generalized in [6] - where a delegated portfolio management problem in continuous time was analyzed. In his model, the Agent

observes prices (a geometric brownian motion) while the Principal observes prices and the fluctuations in wealth resulting from the Agent's investment strategy; investment decisions are unobservable to the Principal and known only to the Agent. Under the assumption of exponential utilities the contracting problem was solved by means of a HJB approach, finding that the optimal contract is of the form "cash plus a convex combination of the generated wealth and a benchmark portfolio"; derivatives are not part of the optimal contract.

Our goal is to clarify the mathematical structure of the optimal contracting problem and to analyze under which conditions on the Principal's and Agent's preferences the main structure-of-equilibrium-contract results in [20] carry over to more general probabilistic settings. To this end, we consider a portfolio delegation model in discrete time, retain the assumptions on the contract space and information structure but allow for rather general utility functions and price dynamics. Our main assumptions on preferences are time-consistency and translation invariance; such preferences have been extensively studied in the mathematical finance literature; see [1, 2, 8, 9, 15]. With our choice of preferences we prove that the problem of dynamic contract design can be reduced to a series of one-period conditional optimization problems of risk-sharing type under constraints and optimal contracts can be computed by backwards induction. To do so, we employ the usual approach of viewing the Agent's continuation utility at any point in time as the Principal's decision variable, with the Principal's decisions being restricted by an incentive compatibility constraint. To the best of our knowledge this argument was first put forward in [24] and later in [23].

Our approach of reducing the dynamic contracting problem to a series of conditional one-period problems is similar to the one employed in [8] where a model of equilibrium pricing in incomplete markets was analyzed. The optimization therein is simpler, though, as the exchange of risk takes place through linear subspaces spanned by the tradable assets which is not the case in our model. Our optimization problems can be viewed as conditional extensions of the ones analyzed in e.g. [1, 5, 8] where the exchange of risk takes place through (conditional) $L^p$ spaces. Conditional analysis - see e.g. [15, 14, 7] - provides a framework to tackle conditional optimization problems, at the same time avoiding technical measurable selection arguments.

Our conditional one-period optimization problems are not convex a-priori, due to incentive compatibility constraints. However, with our choice of contracts the Principal's and Agent's problems can be merged into unconstrained ones, which if solvable yield an optimal contract. In economic terms, the reduction to unconstrained problems means that the first-best solution is implementable under moral hazard if it exists: the contract that one obtains is the same that one would obtain if the Principal and Agent had the same information and had to share the gains and losses from trading between themselves so as to maximize aggregate utility.

The intuition is that in computing an optimal contract the Principal computes the Agent's optimal actions as function of stock prices. This resolves the asymmetry of information and leads eventually to our main result that the optimal contract - if it exists - is of the form "cash plus a convex combination of the generated wealth and a benchmark portfolio plus a path-dependent derivative on the stock price process". In particular, under an optimal contract the Principal fully surrenders to the Agent the wealth generated by trading in exchange for a benchmark portfolio plus a (generally non-replicable) derivative.

As in [20] the benchmark portfolio is related to the optimal (at-equilibrium) effort of the Agent. Unlike in [20] derivatives are generally part of optimal compensation schemes. Derivatives are *not* needed in Markovian models under a predictable representation property (PRP)[1]. The latter includes a discrete-time version of the model in [20] as well as many of the standard dynamic risk sharing problems under symmetric

---

[1]Loosely speaking the predictable representation property states that uncertainty is spanned by finitely many random factors.

information as special cases. It is only in this (restricted) setting that we can prove in Proposition 4.8 that the structure-of-contracts results in [20] carry over to more general preferences as long as the Agent's and the Principal's preferences originate from a common base preference functional (e.g. exponential utilities).

The main challenge is then to solve the unconstrained optimization problems. The approach we follow is to prove that the set of potential optimizers is bounded in a suitable sense. In the greatest generality we work in the conditional version of $L^1$ spaces and with conditional utility functions enjoying a certain sequential upper-semicontinuity on balls, which in particular yields a variational representation of the preference functionals in the spirit of [13, 16, 19]. The transit from boundedness to optimality uses a form of the usual Komlos argument. With this we fully solve in Theorem 2.21 the PA problem for a class of Optimized Certainty Equivalent (OCE) utilities including Average Value at Risk, and bounded prices.

In a Markovian framework under PRP our static conditional problems reduce to deterministic ones in Euclidean spaces. For such setting we find for general OCEs the optimal contract by the Lagrange multiplier method. Under PRP the solution to our contracting problem can also be obtained in terms of the solution to a coupled system of backward stochastic difference equations. As in [8] the benefit of having a discrete model is that such systems can be solved by backwards induction, whereas the continuous-time equivalent is usually intractable. This, and the fact that continuous-time models are unlikely to yield additional insights into the structure of optimal contracts over discrete-time models, motivates our discrete-time framework.

The remainder of this paper is structured as follows. In Section 2 we introduce the modeling framework including the preferences and the contract space. We show how the dynamic problem reduces to a sequence of static ones and present our main results along with examples (OCEs) for which these results can be applied. In Section 3 we prove general attainability results for the Agent's and Principal's problem. In Section 4 we specialize our analysis assuming Markovianity and PRP, which allows us to obtain the optimal contracts explicitly. In Appendix A we survey existing and prove new conditional analysis results which we need throughout this work. Appendices B and C prove results on OCEs and one of the main results of this paper, respectively.

*Notation.* We take the convention that vectors are regarded as column ones. The transpose of a vector $x$ is denoted $x'$ and unless necessary to do otherwise the inner product of two vectors $x, y$ is denoted $xy$. As usual, $(\cdot)_+$ and $(\cdot)_-$ denotes taking positive and negative parts.

## 2. THE MODEL AND MAIN RESULTS

We consider a discrete time model with time grid $\{0, 1, ..., T\}$ for some deterministic terminal time $T < \infty$. Uncertainty is modelled by a probability space $(\Omega, \mathcal{F}, \mathbb{P})$. The probability space carries an $N$-dimensional, strictly positive, discounted stock price process $P = \{P_t\}$ whose filtration we denote by $\mathcal{F} = \{\mathcal{F}_t\}$. We assume throughout that $\mathbb{E}[P_{t+1}|\mathcal{F}_t]$ is finite[2]. We put $\Delta P_{t+1} := P_{t+1} - P_t$ and $\Delta \tilde{P}_{t+1} := diag(P_t)^{-1}\Delta P_{t+1}$, where $diag(\cdot)$ denotes the diagonal matrix associated with the vector in its argument. The same notation applies for other processes different than $P$. We write $P_{0:t}$ to denote the path of the price process from time 0 to $t$. For a $\sigma$-algebra $\mathcal{G}$ we denote by $L^0(\mathcal{G})$ the set of real-valued $\mathcal{G}$-measurable functions. $\underline{L}^0(\mathcal{G})$ and $\overline{L^0}(\mathcal{G})$ denote the set of $\mathcal{G}$-measurable functions taking values in $\mathbb{R} \cup \{-\infty\}$, respectively $\mathbb{R} \cup \{+\infty\}$.

### 2.1. **Effort levels and wealth dynamics.**

At each time $t \in \mathbb{T} := \{0, 1, ..., T-1\}$ the Agent (he) chooses an N-dimensional $\mathcal{F}_t$-measurable random variable $A_t$ that we call *effort level*, in line with the Principal-Agent literature. For the delegated portfolio

---

[2]All equalities and inequalities are to be understood in the $\mathbb{P}$-a.s. sense.

optimization application we have in mind the vector $A_t$ stands for the dollar amount invested in each asset. The cost associated with choosing $A_t$ is given by $c_t(A_t)$. We make the following standing assumption:

**Assumption 2.1.** The cost functions $c_t(\cdot) : \mathbb{R}^N \to \mathbb{R}$ are strictly convex for each $t \in \mathbb{T}$.

Effort levels are known only to the Agent. The wealth at time $t \in \mathbb{T}$ associated with a sequence of effort levels $A = \{A_t\}$ is given by[3]:

$$W_t^A = W_0 + \Delta W_1^A + \cdots + \Delta W_t^A = W_0 + \sum_{s<t} A_s \Delta \tilde{P}_{s+1}. \tag{1}$$

The Principal (she) observes progressively stock prices and wealth levels and offers the Agent a contract based on her available information. Following [20] a contract will consist of a linear combination of a payment contingent on the evolution of the price process and a reward depending linearly on the wealth increments. This includes replicable derivatives on the terminal wealth.

## 2.2. **Preferences.**

Payments are evaluated according to a family of time-consistent and translation invariant utility functions. To connect with the existing literature we first define our preference functionals on spaces of almost surely finite random variables ("general framework"). Subsequently, we introduce an additional conditional integrability and continuity condition ("conditional $L^1$ framework") from which we infer a variational representation of our preference functionals. While time-consistency and translation invariance allows us to reduce the dynamic contracting problem to a series of conditional one-step ones, the variational representation allows us to state sufficient conditions under which the Principal's and the Agent's conditional optimal one-step payments and actions exist at any time.

2.2.1. *General framework.* To introduce our preference functionals we denote by $\mathcal{F}^A$, for a given choice of effort level $A$, the filtration generated by the pair of processes $(P, W^A)$. For the Agent the filtrations $\mathcal{F}$ and $\mathcal{F}^A$ coincide; for the Principal they differ unless she knows the Agent's actions[4]. The respective preferences are then encoded by a family of utility functionals:

$$U_t^a : L^0(\mathcal{F}_T) \to \underline{L}^0(\mathcal{F}_t) \quad \text{and} \quad U_t^p : L^0(\mathcal{F}_T^A) \to \underline{L}^0(\mathcal{F}_t^A) \quad (t \in \mathbb{T}). \tag{2}$$

We use the notation $U^a$ and $U^p$ when referring to the Agent's and Principal's preferences. For a filtration $\{\mathcal{G}_t\}$ and a family $U := \{U_t\}$ of utility functionals $U_t : L^0(\mathcal{G}_T) \mapsto \underline{L}^0(\mathcal{G}_t)$ we say that $U$ is:

- *normalized if* $U_t(0) = 0$,
- *proper if* there exists $X' \in L^0(\mathcal{G}_T)$ s.t. $U_t(X') > -\infty$ and $U_t(X) < \infty$ for all $X \in L^0(\mathcal{G}_T)$
- *monotone if* $U_t(X) \geq U_t(Y)$ whenever $X, Y \in L^0(\mathcal{G}_T)$ and $X \geq Y$
- $\mathcal{F}_t$-*conditionally concave if* $U_t(\lambda X + (1-\lambda)Y) \geq \lambda U_t(X) + (1-\lambda)U_t(Y)$ whenever $\lambda \in L^0(\mathcal{G}_t) \cap [0,1]$ and $X, Y \in L^0(\mathcal{G}_T)$,
- $\mathcal{F}_t$-*translation invariant if* $U_t(X + Y) = U_t(X) + Y$ whenever $X \in L^0(\mathcal{G}_T)$ and $Y \in L^0(\mathcal{G}_t)$,
- *time consistent if* $U_{t+1}(X) \geq U_{t+1}(Y)$ implies $U_t(X) \geq U_t(Y)$,

for all $t \in \mathbb{T}$. We shall refer to these axioms as the *usual conditions/assumptions* and denote by

$$dom(U_t) := \{X \in L^0(\mathcal{G}_T) : U_t(X) \in L^0(\mathcal{G}_t)\},$$

---

[3]For simplicity we assume zero interest rate

[4]The fact that the Principal observes price and wealth dynamics does not necessarily mean that she can observe directly Agent's decisions. For *optimal* contracts the Principal will indeed know Agent's decisions as function of prices. This is not true "off equilibrium", though. Hence we need to distinguish Agent's and Principal's information at this point.

the domain of $U_t$. For a detailed discussion of the usual conditions along with their implications for utility optimization and equilibrium pricing we refer to [8] and references therein. For instance, it is well-known that they imply the *tower property*, stating that $U_t(X) = U_t(U_{t+1}(X))$ whenever $X \in dom(U_{t+1})$, as well as the *local property*, stating that $\mathbb{1}_A U_t(X) = \mathbb{1}_A U_t(Y)$ whenever $X, Y \in L^0(\mathcal{G}_T)$, $A \in \mathcal{G}_t$ and $\mathbb{1}_A X = \mathbb{1}_A Y$.

We assume throughout that $U_t^a$ and $U_t^p$ satisfy the usual conditions w.r.t. the respective filtrations. They are satisfied for a wide class of preferences as illustrated by the following examples.

**Example 2.2 (Entropic utilites).** Given a constant $\gamma > 0$ the entropic family given by

$$U_t(X) = -\frac{1}{\gamma} \log \mathbb{E}\left[\exp\left(-\gamma X\right) | \mathcal{G}_t\right],$$

evidently satisfies the usual conditions.

**Example 2.3 (Pasting).** Starting from one-step utilities defined for bounded random variables, a family satisfying the usual conditions can be built over $L^0$ as follows [8, Example 2]. For each $t$, let $\tilde{U}_t : L^\infty(\mathcal{G}_{t+1}) \to L^\infty(\mathcal{G}_t)$ be a normalized and $\mathcal{G}_t$−translation invariant functional, for which the extensions

$$(3) \qquad\qquad X \mapsto \lim_{n \to +\infty} \lim_{m \to -\infty} \tilde{U}_t([X \wedge n] \vee m),$$

again denoted $\tilde{U}_t$, are well defined between $\underline{L}(\mathcal{G}_{t+1})$ and $\underline{L}(\mathcal{G}_t)$. It is not difficult to see that the *pasting* $U_t(X) := \tilde{U}_t \circ \tilde{U}_{t+1} \circ \cdots \circ \tilde{U}_T(X)$ forms a time consistent and translation invariant family.

**Example 2.4 (Optimized Certainty Equivalents).** Consider $H_t(\cdot)$ a convex, closed and increasing function satisfying $H_t^*(1) := \sup_s[s - H_t(s)] = 0$, and define the one-step functionals

$$\tilde{U}_t : X \in L^\infty(\mathcal{G}_{t+1}) \mapsto \operatorname*{ess\,sup}_{s \in \mathbb{R}}\{s - \mathbb{E}[H_t(s - X)|\mathcal{G}_t]\},$$

which are then normalized, translation-invariant and monotone. Such a family is called *Optimized Certainty Equivalent* (OCE) in the literature; see [2]. The entropic utility of Example 2.2 corresponds to $H(l) = \gamma^{-1} \exp(\gamma l - 1)$. Lemma B.1 shows that if $1 \in int(dom(H_t^*))$ and $H_t$ is bounded from below, the extensions (3) are well-defined. Hence we obtain a family satisfying the usual conditions by pasting; see Remark B.2 as well. This fills a minor gap in [8].

**Example 2.5 (Tail-value-at-risk utility).** By Lemma B.1, a family satisfying all the requirements of Example 2.4 is given by the so-called Tail-value-at-risk (TVAR) utilities, defined for each $\lambda \in (0, 1)$ by

$$(4) \qquad\qquad \tilde{U}_t(X) = \operatorname*{ess\,sup}_s \left\{s - \lambda^{-1}\mathbb{E}([s - X]_+ | \mathcal{G}_t)\right\}.$$

TVAR (or Average-Value-at-Risk) was characterized in [21] and later extensively analyzed in the mathematical finance literature. The representation (4) is more convenient for us than the equivalent:

$$\tilde{U}_t(X) = -\frac{1}{\lambda} \int_{1-\lambda}^1 V@R_\alpha(-X|\mathcal{G}_t)d\alpha.$$

2.2.2. *Conditional $L^1$ framework.* We are now going to introduce additional conditional integrability and continuity conditions on our preference functionals (we refer to Appendix A for more details). We define for two sigma-algebras $\mathcal{G} \subset \tilde{\mathcal{G}}$ the *conditional $L^1$ space*

$$L_{\mathcal{G}}^1(\tilde{\mathcal{G}}) := \left\{Z \in L^0(\tilde{\mathcal{G}}) : \mathbb{E}[|Z||\mathcal{G}] \in L^0(\mathcal{G})\right\}.$$

For $p < \infty$ the $L^p$ variant thereof is evident and we remark that $L_{\mathcal{G}}^1(\tilde{\mathcal{G}}) = L^0(\mathcal{G})L^1(\tilde{\mathcal{G}})$ as sets. Call also

$$L_{\mathcal{G}}^\infty(\tilde{\mathcal{G}}) := \{Z \in L^0(\tilde{\mathcal{G}}) : |Z| \le Y, \text{ for some } Y \in L^0(\mathcal{G})\}.$$

The following continuity property can be viewed as a Fatou property in our conditional framework.

**Definition 2.6.** For $p \in [1, \infty)$ a functional $U : L_{\mathcal{G}}^p(\tilde{\mathcal{G}}) \to \underline{L}^0(\mathcal{G})$ is called $L^0 - L^p$ upper semicontinuous if for each sequence $\{X_n\}_n$ bounded in $L_{\mathcal{G}}^p(\tilde{\mathcal{G}})$ (i.e. $\sup_n \mathbb{E}[|X_n|^p|\mathcal{G}] \in L^0(\mathcal{G})$) such that $X_n \to X$ a.s. it holds that $\limsup U(X_n) \leq U(X)$. We use this terminology even if $U$ is defined in a larger set than $L_{\mathcal{G}}^p(\tilde{\mathcal{G}})$.

We state now a standing assumption on the preferences.

**Assumption 2.7.** Let $U$ stand for $U^a$ or $U^p$ and $\mathcal{G}$ for $\mathcal{F}$ or $\mathcal{F}^A$, respectively. Then $U$ satisfies the usual conditions with respect to $\mathcal{G}$. Moreover, $U_t$ is $L^0 - L^1$ upper semicontinuous for each $t$ and

$$\{X_- : X \in dom(U_t)\} \subset L_{\mathcal{G}_t}^1(\mathcal{G}_T).$$

The following representation result is an immediate consequence of Proposition A.9. It will be used below to prove that the Agent's one-step optimization problems have a solution.

**Proposition 2.8.** Under Assumption 2.7 the following variational representations hold:

$$(5) \qquad U_t^p(X) = \operatorname*{ess\,inf}_{Z \in \mathcal{W}_t^A} \left\{ \mathbb{E}[ZX|\mathcal{F}_t^A] + \alpha_t^p(Z) \right\} \text{ for } X \in L_{\mathcal{F}_t^A}^1(\mathcal{F}_{t+1}^A),$$

$$(6) \qquad U_t^a(X) = \operatorname*{ess\,inf}_{Z \in \mathcal{W}_t} \left\{ \mathbb{E}[ZX|\mathcal{F}_t] + \alpha_t^a(Z) \right\} \quad \text{for } X \in L_{\mathcal{F}_t}^1(\mathcal{F}_{t+1}),$$

where

$$\alpha_t^p : L_{\mathcal{F}_t^A}^\infty(\mathcal{F}_{t+1}^A) \to \overline{L}^0(\mathcal{F}_t^A) \quad \alpha_t^a : L_{\mathcal{F}_t}^\infty(\mathcal{F}_{t+1}) \to \overline{L}^0(\mathcal{F}_t),$$

are the respective conjugates of the utility functionals $U_t^p$ and $U_t^a$, and

$$\mathcal{W}_t^A := \left\{ Z \in L_{\mathcal{F}_t^A}^\infty(\mathcal{F}_{t+1}^A) : Z \geq 0, \mathbb{E}[Z|\mathcal{F}_t^A] = 1 \right\} \quad \mathcal{W}_t := \left\{ Z \in L_{\mathcal{F}_t}^\infty(\mathcal{F}_{t+1}) : Z \geq 0, \mathbb{E}[Z|\mathcal{F}_t] = 1 \right\}.$$

The next example shows that entropic families and pastings of many OCE, such as the TVAR families, fulfilll Assumption 2.7 along with the usual conditions. These are hence the canonical utilities to which our main results in Section 2.4 apply.

**Example 2.9.** The entropic families of Example 2.2 clearly fulfilll Assumptions 2.7. In Lemma B.1 we prove that the TVAR families of Example 2.5, and more generally the OCE families of Example 2.4 for which $1 \in int(dom(H_t^*))$ and $H_t$ is bounded from below, fulfilll Assumptions 2.7 after pasting. In Remark 4.5 we will justify that if the Predictable Representation Property holds, then any OCE satisfies Assumption 2.7.

The variational representation of preferences yields a convenient way to define preference functionals that satisfy Assumption 2.7 by specifying families of "conditionally acceptable models" for both parties.

**Example 2.10.** For a filtration $\{\mathcal{G}_t\}$ let $\mathcal{A}_t \subset L_{\mathcal{G}_t}^\infty(\mathcal{G}_{t+1})$ be a convex set (of conditionally acceptable models) and let $\chi_{\mathcal{A}_t}$ be the associated convex indicator function. Then, the preference functional defined by

$$U_t(X) := \operatorname*{ess\,inf}_{Z \in \mathcal{W}_t} \left\{ \mathbb{E}[ZX|\mathcal{G}_t] - \chi_{\mathcal{A}_t}(X) \right\}$$

satisfies Assumption 2.7 after pasting where $\mathcal{W}_t := \left\{ Z \in L_{\mathcal{G}_t}^\infty(\mathcal{G}_{t+1}) : Z \geq 0, \mathbb{E}[Z|\mathcal{G}_t] = 1 \right\}$

### 2.3. Contracts and optimal actions.

The simplest contracts the Principal may offer the Agent consist of a fixed $\mathcal{F}_T$-measurable (lump-sum) payment $\Theta$, which we may interpret as a financial derivative contingent only on the path of the price process, plus a constant $\beta$ times $W_T^A$. Such contracts (or more exactly, menus of payments) take the form:

$$\bar{S} = \left\{ A \mapsto \bar{S}(A) := \Theta(P_{0:T}) + \beta W_T^A \right\}.$$

Because the Principal observes the wealth and price processes progressively, we shall actually consider a wider family of contracts of the form:

$$S = \left\{ A \mapsto S(A) := \Theta(P_{0:T}) + \sum_{t<T} \beta_t \Delta W_{t+1}^A \right\},$$

where $\beta_t \in L^0(\mathcal{F}_t^A)$ and $\Theta$ is as before, which make better use of her available information. This contract space is rather large and contains replicable path-dependent derivatives on the wealth process. However, as in [20], we shall find that an optimal incentive-compatible contract is indeed of the form $\bar{S}$. This is a consequence of our implicit modeling assumption that the Principal does not seek to infer anything about $A$ from observing $P$ and $W^A$, which we may justify as it being too expensive or time-consuming for the Principal.

We will conveniently refer to a contract as $S$, $(\Theta, \beta)$ or $(\Theta, \{\beta_t\})$ depending on the context and denote by $R \in \mathbb{R}$ the Agent's reservation utility, i.e. the least utility the Agent demands in order to commit to a contract

**Definition 2.11.** A contract $(\Theta, \{\beta_t\})$ is *individually rational* if the optimal utility the Agent can obtain at time 0 from it is at least $R$.

In the sequel we show how to obtain recursive representations of the Agent's and the Principal's utilities and how to reduce the problem of optimal dynamic contract design to a sequence of static problems.

2.3.1. *Agent's problem.* Let us assume that the Agent chooses an effort level $A$ when presented with a contract $S(\cdot)$. His total cost of effort is then $C(A) := \sum_{t=0}^{T-1} c_t(A_t)$ and his utility seen from time $t$ is $U_t^a(S(A) - C(A))$. Using translation invariance we compute:

$$U_t^a \left( S(A) - \sum_t c_t(A_t) \right) = U_t^a \left( \Theta(P_{0:T}) + \sum_{s \geq t} \left\{ \beta_s \Delta W_{s+1}^A - c_{s+1}(A_{s+1}) \right\} \right)$$

(7)
$$- c_t(A_t) + \sum_{s<t} \left\{ \beta_s \Delta W_{s+1}^A - c_s(A_s) \right\}.$$

This shows that the Agent's optimization problem of finding the best effort level $A$ given a contract $S(\cdot)$ reduces to the following recursion (we omit for simplicity the dependence of $H$ in $S$):

$$H_T = \Theta(P_{0:T})$$
(8)
$$H_t = \operatorname*{ess\,sup}_{A \in L^0(\mathcal{F}_t)^N} \left\{ U_t^a \left( H_{t+1} + \beta_t A \Delta \tilde{P}_{t+1} \right) - c_t(A) \right\}.$$

*Remark* 2.12. The preceding analysis shows that $H_t$ has the interpretation of being the maximal utility the Agent can get, from time $t$ onwards. Since adding an $\mathcal{F}_t-$measurable term to $\Theta$ translates additively into $H_t$ and preserves optimality of effort levels, we see that the individual rationality condition binds $(H_0 = R)$ for any contract that is optimal for the Principal.

**Definition 2.13.** A contract $(\Theta, \{\beta_t\})$ is called *incentive-compatible* if the essential suprema in (8) are attained for each $t \in \mathbb{T}$.

2.3.2. *Principal's problem.* The Principal's problem is to design an optimal incentive compatible and individually rational contract. To that end, suppose again that the Agent has chosen $A$ when presented

with a contract $S(\cdot)$, and that the Principal knows this. Her utility seen from time $t$ is then:

$$U_t^p \left( W_T^A - \Theta - \sum_{s<T} \beta_s \Delta W_{s+1}^A \right)$$

$$= W_0^A - H_t + \sum_{s<t}(1-\beta_s)A_s\Delta\tilde{P}_{s+1} + U_t^p \left( \sum_{s\geq t} \left[ (1-\beta_s)A_s\Delta\tilde{P}_{s+1} - \Delta H_{s+1} \right] \right),$$

where the identity $\Theta = H_t + \sum_{s\geq t}\Delta H_{s+1}$ and translation invariance was used. If we denote by $h_t(A,\beta)$ her utility from future income, then time consistency along with translation invariance yields:

$$
\begin{aligned}
(9) \quad h_t(A,\beta) &:= U_t^p \left( \sum_{s\geq t} \left[ (1-\beta_s)A_s\Delta\tilde{P}_{s+1} - \Delta H_{s+1} \right] \right) \\
&= U_t^a \left( H_{t+1} + \beta_t A_t \Delta\tilde{P}_{t+1} \right) - c_t(A_t) + U_t^p \left( h_{t+1}(A,\beta) + (1-\beta_t)A_t\Delta\tilde{P}_{t+1} - H_{t+1} \right).
\end{aligned}
$$

Performing the change of variables

$$(10) \qquad\qquad \Gamma_{t+1} := \beta_t A_t \Delta\tilde{P}_{t+1} + H_{t+1} \in L^0(\mathcal{F}_{t+1}),$$

and writing $h_t(A,\Gamma)$ instead of $h_t(A,\beta)$ we arrive at:

$$(11) \qquad h_t(A,\Gamma) = U_t^a(\Gamma_{t+1}) - c_t(A_t) + U_t^p \left( h_{t+1}(A,\Gamma) + A_t\Delta\tilde{P}_{t+1} - \Gamma_{t+1} \right).$$

If $(\Theta,\{\beta_t\})$ is incentive compatible, then unique optimal effort levels for the Agent exist, due to our concavity assumptions on his utility and cost function. For every time $t \in \mathbb{T}$ we may thus construct the random variable $\Gamma_{t+1}$, and $A_t$ will attain the essential supremum:

$$\operatorname*{ess\,sup}_a \left[ U_t^a \left( \Gamma_{t+1} + \beta_t[a - A_t]\Delta\tilde{P}_{t+1} \right) - c_t(a) \right].$$

We say that $(\{A\},\{\Gamma\})$ is *incentive-compatible* whenever for every $t$ this $A_t$ attains this supremum. In terms of the set

$$\mathbb{C}_t(\beta) := \left\{ \begin{array}{c} (A,\Gamma) \in [L^0(\mathcal{F}_t)]^N \times L^0(\mathcal{F}_{t+1}) \text{ s.t. for every } \bar{A} \in [L^0(\mathcal{F}_t)]^N : \\ U_t^a(\Gamma) - c_t(A) \geq U_t^a \left( \Gamma + \beta[\bar{A} - A]\Delta\tilde{P}_{t+1} \right) - c_t(\bar{A}) \end{array} \right\},$$

incentive compatibility amounts to $(A_t,\Gamma_{t+1}) \in \mathbb{C}_t(\beta_t)$ for every $t \in \mathbb{T}$. In particular, we can introduce the following recursion for the Principal's future optimal wealth:

$$h_T = 0,$$

$$(12) \qquad h_t = \operatorname*{ess\,sup}_{\substack{(\beta,A,\Gamma) \\ (A,\Gamma)\in\mathbb{C}_t(\beta)}} U_t^a(\Gamma) - c_t(A) + U_t^p \left( h_{t+1} + A\Delta\tilde{P}_{t+1} - \Gamma \right).$$

*Remark* 2.14. We arrived at the well-known result that in constructing an optimal contract the Principal should consider the Agent's continuation utility as a decision variable of hers. This also resolves the issue of information asymmetry: assuming that the Principal knows the mappings $A_t$ as functions of $\{P_s\}_{s\leq t}$ for each $t$ implies that all the random variables in (9) and (11) become price-adapted.[5] If optimal efforts are not unique, then one has to specify which effort levels $\{A_t\}$ (the Principal recommends) the Agent implements in order to carry out the above recursion. This is why in the PA literature one often calls such effort levels *recommended effort levels* and the triple $(\Theta,\{\beta_t\},\{A_t\})$ incentive compatible.

---

[5]We emphasize, again, that this is a consequence of the assumption that the Principal is not trying to learn/infer something from the Agent's actions.

2.4. **Main results.**

In this section we summarize the main results of our paper. We start with the following theorem that makes our formal derivations of the Agent's and Principal's optimal wealth precise. It states that if the Principal's and the Agent's conditional one-step optimization problems have solutions, then the dynamic contracting problem has a solution that can be obtained out of these. The proof is given in Appendix C.

**Theorem 2.15.** Assume that the recursions (8) and (12) admit a solution with the essential suprema attained at each time $t$. Then the Principal's optimal utility at time $t = 0$ equals $W_0 - R + h_0$. Further, calling $(\beta_t, A_t, \Gamma_{t+1})_{t \in \mathbb{T}}$ the maximizers attaining $h$ in (12), and defining

$$\Theta = \Theta(P_{0:T}) := \sum_{0 \leq t < T} \left[ \Gamma_{t+1} - U_t^a(\Gamma_{t+1}) + c_t(A_t) \right],$$

the contract

$$S = \left\{ \bar{A} \mapsto R + \Theta(P_{0:T}) + \sum \beta_t \left[ \Delta W_{t+1}^{\bar{A}} - A_t \Delta \tilde{P}_{t+1} \right] \right\},$$

is optimal for the Principal, among those satisfying incentive compatibility and individual rationality constraints. The associated optimal effort for the Agent is $A$ and his optimal wealth will be $R$.

We now define an auxiliary unconstrained version of the optimization problem in (12), and prove that if such a problem is well-posed, it yields a at time $t$ a solution to the original one-step problem, and the corresponding $\beta_t = 1$ is optimal. This opens the way to our main result, Theorem 2.21. The technical importance of this is that we may dispense with the non-convex sets $\mathbb{C}_t$, making the incentive-compatibility constraint much more tractable. Economically, this indicates that the first-best solution is optimal if it exists: at any point in time $t \in \mathbb{T}$ both parties share the "aggregate endowment" given by the Principal's utility from future income $h_{t+1}$ plus gains from trading $A \Delta \tilde{P}_{t+1}$ so as to maximize aggregate utility.

**Proposition 2.16.** Assume that the following problem is finite and attainable:

$$(13) \qquad \Sigma := \operatorname*{ess\,sup}_{(A,\Gamma) \in [L^0(\mathcal{F}_t)]^N \times L^0(\mathcal{F}_{t+1})} U_t^a(\Gamma) - c_t(A) + U_t^p \left( h_{t+1} + A \Delta \tilde{P}_{t+1} - \Gamma \right).$$

Then any maximizer $\left( \hat{A}, \hat{\Gamma} \right)$ belongs to the set $\mathbb{C}_t(1)$ and therefore

$$\Sigma = \operatorname*{ess\,sup}_{\substack{(\beta, A, \Gamma) \\ (A,\Gamma) \in \mathbb{C}_t(\beta)}} U_t^a(\Gamma) - c_t(A) + U_t^p \left( h_{t+1} + A \Delta \tilde{P}_{t+1} - \Gamma \right).$$

*Proof.* Let $(\hat{A}, \hat{\Gamma})$ be a maximizer for (13). For arbitrary $A$, define $\Gamma = \hat{\Gamma} + (A - \hat{A}) \Delta \tilde{P}$. Plugging in that $(\hat{A}, \hat{\Gamma})$ is better than $(A, \Gamma)$ for (13), we see that the terms involving $U^p$ cancel out and so:

$$(14) \qquad U_t^a(\hat{\Gamma}) - c_t(\hat{A}) \geq U_t^a(\hat{\Gamma} + (A - \hat{A}) \Delta \tilde{P}) - c_t(A).$$

This means that $\left( \hat{A}, \hat{\Gamma} \right) \in \mathbb{C}_t(1)$ so the values of the constrained and unconstrained problems coincide. $\square$

*Remark* 2.17. The previous proof crucially relies on the fact that contracts are linear in wealth increments. Indeed by varying $\hat{\Gamma}$ in directions of the form $(A - \hat{A}) \Delta \tilde{P}$ and by linearity of contracts the term in the objective function involving Principal's utility cancel out, making it possible to compare the values of Agent's utilities.

In Section 3 we shall, therefore, turn our attention to the question of attainability of the unconstrained problem. For the reader's convenience we state in this section our main results therein and show how they apply to specific classes of examples. The proof of the following result is given in Section 3.2. The technical conditions will be easily satisfied by the utility functionals listed in Example 2.9.

**Theorem 2.18.** Suppose at time $t \in \mathbb{T}$ that

$$K_t^p := \operatorname*{ess\,sup}_{Z \in \mathcal{W}_t \cap [1-\epsilon, 1+\epsilon]} \alpha_t^p(Z) \in L^0(\mathcal{F}_t) \quad \text{and} \quad K_t^a := \operatorname*{ess\,sup}_{Z \in \mathcal{W}_t \cap [1-\epsilon, 1+\epsilon]} \alpha_t^a(Z) \in L^0(\mathcal{F}_t),$$

for some $\epsilon \in L^0(\mathcal{F}_t) \cap (0, 1]$. Then, if $h_{t+1} \in dom(U_t^p)$ and $\lim_{|a| \to \infty} \frac{c_t(a)}{|a|} = +\infty$, the random variable $\Sigma$ defined in (13) belongs to $L^0(\mathcal{F}_t)$, satisfies

$$\Sigma = \operatorname*{ess\,sup}_{(A, \Gamma) \in [L^0(\mathcal{F}_t)]^N \times L^1_{\mathcal{F}_t}(\mathcal{F}_{t+1})} U_t^a(\Gamma) - c_t(A) + U_t^p\left(h_{t+1} + A\Delta\tilde{P}_{t+1} - \Gamma\right),$$

and the essential supremum is attained. In particular $\beta_t = 1$ is optimal at time $t \in \mathbb{T}$.

It is well-known that if the utility functionals originate from a common base functional, more explicit treatments of equilibrium/risk-sharing problems become available (as in [1, 5, 8]). In the same spirit we have the following result, stating that in that case the Principal and the Agent share the "aggregate endowment" according to their risk attitudes. The proof is given in Section 3.2.

**Theorem 2.19** (Base Preferences). Suppose that there exists non-negative numbers $\gamma^a, \gamma^p$ and base preference functionals $\{U_t\}$ such that

$$U_t^l(\cdot) := \frac{1}{\gamma^l} U_t\left(\gamma^l \cdot\right) \quad (l = a, p).$$

Further assume that

$$\frac{\gamma^a \gamma^p}{\gamma^a + \gamma^p} h_{t+1} \in dom(U_t) \quad \text{and} \quad \lim_{|a| \to \infty} \frac{c_t(a)}{|a|} = +\infty.$$

Then Principal's one-step problem (at time $t$) has as solution:

$$\beta = 1 \quad \text{and} \quad \Gamma^* = \frac{\gamma^p}{\gamma^a + \gamma^p}(h_{t+1} + A^*\Delta\tilde{P}_{t+1}),$$

for the optimal action $A^*$ of the Agent, which attains:

$$\operatorname*{ess\,sup}_{A}\left\{-c_t(A) + \frac{\gamma^a + \gamma^p}{\gamma^a \gamma^p} U_t\left(\frac{\{\gamma^a \gamma^p\}[h_{t+1} + A\Delta\tilde{P}_{t+1}]}{\gamma^a + \gamma^p}\right)\right\}.$$

In light of Theorem 2.15, the two previous results yield a solution to the dynamic problem, as explained in the following proposition.

**Proposition 2.20.** If the assumptions of Theorem 2.18 or Theorem 2.19 hold for every $t \in \mathbb{T}$, then the respective one-step problems have a solution and glueing them together yields a solution for the respective dynamic problems, whereby $\beta_t = 1$ for every $t \in \mathbb{T}$ is optimal.

The proof of the preceding proposition is obvious. In applying this result, several technical conditions need be checked *a-posteriori*. As shown by the following theorem, these conditions are satisfied *a-priori* for entropic and TVAR families and for OCE utilities (Example 2.9) under mild conditions. The proof is given in Appendix B.

**Theorem 2.21.** Suppose that prices are bounded ($0 < p_- \leq P_t^i \leq p_+$ a.s.) and that both $U^a$ and $U^p$ are constructed by pasting of optimized certainty equivalent functionals:

$$X \in L^1_{\mathcal{F}_t}(\mathcal{F}_{t+1}) \mapsto U_t^a(X) = \operatorname*{ess\,sup}_{s \in \mathbb{R}}\{s - \mathbb{E}[H_t^a(s - X)|\mathcal{F}_t]\}$$

$$X \in L^1_{\mathcal{F}_t^A}(\mathcal{F}_{t+1}^A) \mapsto U_t^p(X) = \operatorname*{ess\,sup}_{s \in \mathbb{R}}\{s - \mathbb{E}[H_t^p(s - X)|\mathcal{F}_t^A]\},$$

for which the following conditions hold for each $t$:

- $1 \in int(dom(H_t^a)) \cap int(dom(H_t^p))$,
- $H_t^a$ and $H_t^p$ are lower-bounded.

Finally assume that $\lim_{|a| \to \infty} \frac{c_t(a)}{|a|} = \infty$ for every $t$. Then our dynamic Principal-Agent problem has a solution whereby the Agent keeps the output wealth and the Principal is given a possibly path-dependent derivative.

*Remark* 2.22. In conjunction with Theorem 2.15, the previous result yields the economic interpretation we referred to in the introduction: the optimal contract is of the form "cash plus a path-dependent derivative on the stock price process plus performance w.r.t. a benchmark portfolio". As the derivative may not be replicable, this shows that the structure in [20, Theorem 1] need not hold.

A family of examples where we can provide explicitly the form of an optimal contract, recovering the results of [20] in the continuous case is given in Section 4.1 below. It requires additional notation, though, so we postpone the statement of the result to Section 4.

*Remark* 2.23. For simplicity and ease of exposition we took zero interest rates and $c_t = c_t(A_t)$. The case with non-null interest rates and/or $c(A, W) = \sum_t [c_t(A_t) + \gamma_t W_{t-1}]$ can be solved exactly in the same way, the only difference being that $\beta_t$ will not be constant (but remains deterministic) anymore. The qualitative structure of contracts and their interpretation remain the same however.

## 3. GENERAL ATTAINABILITY RESULTS

We prove in this section the attainability of the Agent's and Principal's one-step problems, and consequently, for the dynamic problem.

### 3.1. **Agent's Problem.**

We start with an abstract conditional optimization problem of which the Agent's one-step optimization problems are special cases. For a given pair of random variables $(X, \beta) \in L^0(\mathcal{F}_{t+1}) \times L^0(\mathcal{F}_t)$, let

$$
\begin{aligned}
G(t, X, \beta) &:= \operatorname*{ess\,sup}_{A \in L^0(\mathcal{F}_t)^N} \left\{ -c_t(A) + U_t^a \left( X + \beta A \Delta \tilde{P}_{t+1} \right) \right\} \\
&=: \operatorname*{ess\,sup}_{A \in L^0(\mathcal{F}_t)^N} g_t(A).
\end{aligned}
$$
(15)

Under the usual conditions $g_t$ is $\mathcal{F}_t-$concave, and hence stable (see Definition A.4). The key to the above optimization problem is to reduce it to an $L^0(\mathcal{F}_t)-$bounded set.

**Lemma 3.1.** Under the following condition, the essential supremum in (15) is attained:

$$
X \in dom(U_t^a) \quad \text{and} \quad \lim_{|a| \to \infty} \frac{c_t(a)}{|a|} = +\infty.
$$

*Proof.* We intend to apply Theorem A.6. Evidently

$$
\operatorname*{ess\,sup}_{A \in [L^0(\mathcal{F})]^N} g_t(A) = \operatorname*{ess\,sup}_{A \in \Lambda} g_t(A),
$$

where $\Lambda = \{A : g_t(A) \geq g_t(0)\}$. The set $\Lambda$ is $L^0-$convex, contains the origin and is $\sigma-$stable. That $\Lambda$ is sequentially closed is an application of Proposition A.9.

For $A \in [L^0(\mathcal{F}_t)]^N$ not identically null we use the variational representation of $U^a$ established in Proposition 2.8 to bound:

$$
\begin{aligned}
g_t(nA) &= U_t^a(X + n\beta A \Delta \tilde{P}) - c_t(nA) \\
&\leq K + \mathbb{E}[ZX_+|\mathcal{F}_t] + n\mathbb{E}[\beta Z A \Delta \tilde{P}|\mathcal{F}_t] - c_t(nA),
\end{aligned}
$$
(16)

where $Z \in \mathcal{W}_t$. Using that $A, \beta$ are $\mathcal{F}_t$-measurable and Cauchy-Schwarz applied pointwise, we bound from above the sum of the last two terms in (16) on the set where $A$ does not vanish by:

$$n|A||\beta||\mathbb{E}[Z\Delta\tilde{P}|\mathcal{F}_t]| - c_t(nA) \leq n|A| \left[ |\beta||\mathbb{E}[Z\Delta\tilde{P}|\mathcal{F}_t]| - \frac{c_t(nA)}{n|A|} \right].$$

Since $Z \in L^\infty_{\mathcal{F}_t}(\mathcal{F}_{t+1})$ and $|\mathbb{E}[\Delta\tilde{P}|\mathcal{F}_t]|$ is a.s. finite by assumption, we see that the majorizing term tends to $-\infty$ on a non-negligible set as $n \to \infty$ and so does $g_t(nA)$. Since $g_t(0) = U^a_t(X) - c_t(0) > -\infty$ by assumption, we get a contradiction, and so Theorem A.3 shows that $\Lambda$ is $L^0(\mathcal{F}_t)$-bounded. Hence Theorem A.6 applies to $\operatorname{ess\,sup}_{A \in \Lambda} g_t(A)$, since the mapping $A \mapsto g_t(A)$ is $L^0$-upper semicontinuous by Proposition A.9. This establishes attainability. $\qquad\square$

The following is an immediate corollary of the previous lemma.

**Corollary 3.2.** Assume that $H_{t+1} \in dom(U^a_t)$ and $\lim_{|a| \to \infty} \frac{c_t(a)}{|a|} = +\infty$. Then the one-step conditional optimization problem of the Agent at time $t$, as in (8), is attained.

### 3.2. **Principal's Problem.**

In this section we prove Theorem 2.18, which sharpens Proposition 2.16. The Principal's problem at time $t$ consists in maximizing

$$V_t(A, \Gamma) := U^a_t(\Gamma) - c_t(A) + U^p_t \left( h_{t+1} + A\Delta\tilde{P}_{t+1} - \Gamma \right).$$

Recall from Remark 2.14 that the Principal's preference functionals $U^p_t$ may and will be considered as mappings from $L^0(\mathcal{F}_T)$ to $L^0(\mathcal{F}_t)$, satisfying the usual assumptions w.r.t. $\mathcal{F}$.

*Proof of Theorem 2.18.* Let us introduce the set

$$\mathcal{S} := \{(A, \Gamma) \in L^0(\mathcal{F}_t)^N \times Q : V(A, \Gamma) \geq V(0, 0)\},$$

where $Q := \left\{ \Gamma \in L^1_{\mathcal{F}_t}(\mathcal{F}_{t+1}) : \mathbb{E}[\Gamma|\mathcal{F}_t] = 0 \right\}$. In maximizing $V$, i.e. in computing $\Sigma$, we may assume that $\Gamma \in L^1_{\mathcal{F}_t}(\mathcal{F}_{t+1})$, since for candidate optima, $\Gamma$ and $-\Gamma$ must be in the domains of $U^a$ and $U^p$ respectively, and by assumption this yields $\Gamma_-, \Gamma_+ \in L^1_{\mathcal{F}_t}(\mathcal{F}_{t+1})$. We may thus further assume that $(A, \Gamma)$ belong to $\mathcal{S}$, since $\mathcal{F}_t$-measurable components of $\Gamma$ cancel out in $V$, i.e.

$$\Sigma = \operatorname*{ess\,sup}_{(A,\Gamma) \in \mathcal{S}} V(A, \Gamma).$$

In a first step, we will show that the set

$$\mathcal{S}^A := \{A \in L^0(\mathcal{F}_t)^N : \text{ there exists } \Gamma \in Q \text{ such that } (A, \Gamma) \in \mathcal{S}\}$$

is $L^0(\mathcal{F}_t)-$bounded. To this end, we first notice that $V(0, 0) = -c(0) + U^p_t(h_{t+1}) \in L^0(\mathcal{F}_t)$. Taking

$$\tilde{Z} \in dom(\alpha^p) \cap dom(\alpha^a) \cap \mathcal{W} \cap L^0(\mathcal{F}_t)$$

(e.g. $\tilde{Z} = 1$) and using the variational representation of the preference functionals, we get:

$$
\begin{aligned}
U^a_t(\Gamma) &\leq \alpha^a(\tilde{Z}) + \tilde{Z}\mathbb{E}[\Gamma|\mathcal{F}_t] \\
U^p_t(h + A\Delta\tilde{P} - \Gamma) &\leq \alpha^p(\tilde{Z}) + \tilde{Z}\mathbb{E}[h|\mathcal{F}_t] - \tilde{Z}\mathbb{E}[\Gamma|\mathcal{F}_t] + \tilde{Z}\mathbb{E}[A\Delta\tilde{P}_{t+1}|\mathcal{F}_t].
\end{aligned}
$$

For $\Gamma \in Q$ the term $\mathbb{E}[\Gamma|\mathcal{F}_t]$ vanishes and hence

$$V(0, 0) \leq \alpha^p(\tilde{Z}) + \alpha^a(\tilde{Z}) + \tilde{Z}\mathbb{E}[h|\mathcal{F}_t] + |A||\tilde{Z}\mathbb{E}[\Delta\tilde{P}_{t+1}|\mathcal{F}_t]| - c(A).$$

Since $\mathcal{S}^A$ is $\sigma$-stable, we can use Lemma A.5 to conclude. Indeed, if $\mathcal{S}^A$ were not $L^0(\mathcal{F}_t)$-bounded, then there exists a non-negligible set $\tilde{\Omega}$ and a sequence $\{A_n\} \subset \mathcal{S}^A$ such that $|A_n| \geq n$ on $\tilde{\Omega}$. Similar arguments

as in the proof of Lemma 3.1 would establish $V(0,0) = -\infty$ on a non-negligible set, contradicting our hypotheses. Thus $\mathcal{S}^A$ must be $L^0(\mathcal{F}_t)$-bounded.

Next, we prove that the set

$$\mathcal{S}^\Gamma := \{\Gamma \in Q : \text{ there exists } A \in L^0(\mathcal{F}_t)^N \text{ such that } (A, \Gamma) \in \mathcal{S}\}$$

is bounded in $L^1_{\mathcal{F}_t}(\mathcal{F}_T)$. Let us chose $\epsilon \in L^0(\mathcal{F}_t) \cap (0, 1]$ as in the statement of this theorem, fix $\Gamma \in \mathcal{S}^\Gamma$ and define

$$Z^a := 1 + \epsilon[\mathbb{1}_{\Gamma \leq 0} - \mathbb{P}(\Gamma \leq 0 | \mathcal{F}_t)] \in L^\infty(\mathcal{F}_T) \cap [1 - \epsilon, 1 + \epsilon]$$
$$Z^p := 1 + \epsilon[\mathbb{1}_{\Gamma > 0} - \mathbb{P}(\Gamma > 0 | \mathcal{F}_t)] \in L^\infty(\mathcal{F}_T) \cap [1 - \epsilon, 1 + \epsilon]$$

Since $\Gamma \in Q$ we see that

$$\mathbb{E}[Z^a \Gamma | \mathcal{F}_t] = -\epsilon \mathbb{E}[(\Gamma)_- | \mathcal{F}_t] \quad \text{and} \quad \mathbb{E}[Z^p \Gamma | \mathcal{F}_t] = \epsilon \mathbb{E}[(\Gamma)_+ | \mathcal{F}_t].$$

Moreover, $\mathbb{E}[Z^a | \mathcal{F}_t] = \mathbb{E}[Z^p | \mathcal{F}_t] = 1$, implying that $Z^{a,p} \in \mathcal{W}_t$ and thus $\alpha^p(Z^p) \leq K^p$ and $\alpha^a(Z^a) \leq K^a$. We hence obtain that

$$U^a(\Gamma) \leq -\epsilon \mathbb{E}[(\Gamma)_- | \mathcal{F}_t] + K^a,$$

$$U^p(h + A\Delta\tilde{P} - \Gamma) \leq \mathbb{E}[Z^p(h + A\Delta\tilde{P}) | \mathcal{F}_t] - \epsilon \mathbb{E}[\Gamma_+ | \mathcal{F}_t] + K^p$$
$$\leq 2\mathbb{E}[|h| | \mathcal{F}_t] + 2|A| \mathbb{E}[|\Delta\tilde{P}| | \mathcal{F}_t] - \epsilon \mathbb{E}[\Gamma_+ | \mathcal{F}_t] + K^p$$
$$\leq N - \epsilon \mathbb{E}[\Gamma_+ | \mathcal{F}_t],$$

for some $N \in L^0(\mathcal{F}_t)$ where the latter inequality follows by assumption and the fact that the effort levels had already been proven to be $L^0(\mathcal{F}_t)$-bounded. Therefore for $(A, \Gamma) \in \mathcal{S}$ we have

$$V(0,0) \leq V(A, \Gamma) \leq N + K^a - \epsilon E[(\Gamma)_- | \mathcal{F}_t] - \epsilon \mathbb{E}[\Gamma_+ | \mathcal{F}_t] - c_t(A) \leq \tilde{K} - \epsilon E[|\Gamma| | \mathcal{F}_t],$$

for some $\tilde{K} \in L^0(\mathcal{F}_t)$. This implies that $\mathcal{S}^\Gamma$ is bounded in $L^1_{\mathcal{F}_t}(\mathcal{F}_T)$ since $\epsilon > 0$ a.s.

Next, we notice that there exists a sequence $(A_n, \Gamma_n) \in \mathcal{S}$ such that $V(A_n, \Gamma_n) \uparrow \Sigma$ since $\mathcal{S}$ is directed upwards. Indeed, if $V(A^i, \Gamma^i) \geq V(0,0)$ for $i = 1, 2$ and if we define $\xi = \{V(A^1, \Gamma^1) \geq V(A^2, \Gamma^2)\}$ and $(A, \Gamma) = (A^1, \Gamma^1)\mathbb{1}_\xi + (A^2, \Gamma^2)\mathbb{1}_{\xi^c}$, then

$$V(A, \Gamma) = \max\{V(A^1, \Gamma^1), V(A^2, \Gamma^2)\} \geq V(0,0),$$

thanks to the terms in $V$ being $\mathcal{F}_t$–stable and $\xi \in \mathcal{F}_t$. By virtue of $\mathcal{S}^A$ being $L^0(\mathcal{F}_t)$-bounded, we can apply the usual Komlos lemma (or Lemma A.7) to the positive and negative parts of each component of the sequence $\{A_n\}_n$ in an iterative, nested way, i.e. taking convex combinations of convex combinations and so forth. On the other hand, the $L^1_{\mathcal{F}_t}(\mathcal{F}_T)$-boundedness of $\mathcal{S}^\Gamma$ implies that the technical condition in Lemma A.7 holds for the positive and negative parts of the sequence $\{\Gamma_n\}_n$, by Jensen's inequality, so we can again take convex combinations of convex combinations. All in all we have found a sequence of non-negative real numbers $\{\lambda_i^n\}$ with $\sum_{i \geq n} \lambda_i^n = 1$, and random variables $\Gamma^* \in L^0(\mathcal{F}_{t+1})$ and $A^* \in L^0(\mathcal{F}_t)^N$ such that $\tilde{\Gamma}_n = \sum_{i \geq n} \lambda_i^n \Gamma_i \to \Gamma^*$ and $\tilde{A}_n = \sum_{i \geq n} \lambda_i^n A_i \to A^*$ a.s. (for each component). Also $(\tilde{A}_n, \tilde{\Gamma}_n) \in \mathcal{S}$ by convexity. Moreover,

$$\Sigma = \lim_n V(A_n, \Gamma_n) = \lim_n \sum_{i \geq n} \lambda_i^n V(A_i, \Gamma_i) \leq \limsup_n V(\tilde{A}_n, \tilde{\Gamma}_n),$$

since (a.s.) convergent sequences of real numbers remain converging under convex combinations of its tails and $V$ is concave.

The cost-term in $V$ is u.s.c. and since $\mathcal{S}^\Gamma$ is $L^1_{\mathcal{F}_t}(\mathcal{F}_T)$-bounded we get for the $U^a$ term in $V$ that $\limsup_n U^a\left(\tilde{\Gamma}_n\right) \leq U^a(\Gamma^*)$. Finally, for the $U^p$ term in $V$, we obtain from the last assertion in Proposition A.9 that $\limsup_n U^p\left(h_{t+1} + \tilde{A}_n \Delta\tilde{P}_{t+1} - \tilde{\Gamma}_n\right) \leq U^p(h_{t+1} + A^* \Delta\tilde{P}_{t+1} - \Gamma^*)$. We thus get that

$\Sigma \leq V(A^*, \Gamma^*)$ and hence we have equality. This shows that $\Sigma < \infty$ since the preference functionals are proper. Finally, by Proposition 2.16 we conclude that $\beta = 1$ is optimal and Principal's one-step problem is attained. $\qquad\square$

We proceed now to the proof of Theorem 2.19.

*Proof of Theorem 2.19.* Let us first fix an effort level $A$ and put $x := h + A\Delta\tilde{P}_{t+1}$ and $\hat{\gamma} := \frac{\gamma^a \gamma^p}{\gamma^a + \gamma^p}$. Concavity of the preference functional yields:

$$\operatorname*{ess\,sup}_{\Gamma} \left\{ U_t^a(\Gamma) + U_t^p(x - \Gamma) \right\} = \operatorname*{ess\,sup}_{\Gamma} \frac{1}{\hat{\gamma}} \left\{ \frac{\hat{\gamma}}{\gamma^a} U_t(\gamma^a \Gamma) + \frac{\hat{\gamma}}{\gamma^p} U_t(\gamma^p[x - \Gamma]) \right\} \leq \frac{1}{\hat{\gamma}} U_t(\hat{\gamma} x).$$

On the other hand, taking $\Gamma^* = \frac{\gamma^p}{\gamma^a + \gamma^p} x$ it follows that $\frac{1}{\gamma^a} U_t(\gamma^a \Gamma^*) + \frac{1}{\gamma^p} U_t(\gamma^p[x - \Gamma^*]) = \frac{1}{\hat{\gamma}} U_t(\hat{\gamma} x)$. Therefore this $\Gamma^*$ attains the essential supremum above. Thus the Principal's problem reduces to:

$$(17) \qquad \operatorname*{ess\,sup}_{\Gamma} \left\{ -c_t(A) + \frac{1}{\hat{\gamma}} U_t \left( \hat{\gamma}[h_{t+1} + A\Delta\tilde{P}_{t+1}] \right) \right\}.$$

If this problem is attained at $A^*$, then the previous argument shows that $\Gamma^* = \frac{\gamma^p}{\gamma^a + \gamma^p}(h_{t+1} + A^*\Delta\tilde{P}_{t+1})$ is optimal. The problem (17) is of the same form as that analyzed in Lemma 3.1, simply replacing $U^a$ by $\frac{1}{\hat{\gamma}} U_t(\hat{\gamma}\cdot)$, calling $X = h$ and taking $\beta = 1$. In particular, we obtain existence of an optimizer $A^*$. Because the one-step unconstrained problem is attained, Proposition 2.16 shows that taking $\beta = 1$, $A^*$ and $\Gamma^*$ at time $t$ yields an optimal one-step decision. $\qquad\square$

*Remark* 3.3. In this article we chose to work in the biggest conditional (loc. convex) space of $L^p$-type, this is, the conditional $L^1$ space. The reason is twofold. On the one hand, had we worked with smaller subspaces, we would have had in principle more tools at hand to prove the attainability of Principal's one-step problems. However, we chose not to limit the scope of utility functionals a priori, in terms of their domains, for which the theory would be applicable to. On the other hand, even acknowledging the fact that our $L^0 - L^1$ upper semicontinuity requirement is not a mild one, the alternative would have been to impose from the outset some sort of "sup-compactness" of our functionals (more precisely, of their convolutions) or again to work with smaller spaces than conditional $L^1$; ideally conditionally reflexive ones. It seems to us that our simple sequential (and rather point-wise) $L^0 - L^1$ upper semicontinuity has the advantage of being a more tractable and less technical requirement than the other, very valid ones.

## 4. Optimal contracting under predictable representation

Up to now our probability space and price process were rather general. In this section we add more structure to the problem in order to obtain more explicit solutions. In particular we fix a volatility matrix $\sigma \in \mathbb{R}^{N,d}$ with linearly independent rows ($d \geq N$), assume that the flow of information is generated by a d-dimensional process $\bar{w} = (w^1, ..., w^d)$ whose evolution is observed by both parties and that the price dynamics follows:

$$(18) \qquad \Delta P_{t+1} = diag(P_t)\left[\mu + \sigma\Delta\bar{w}_{t+1}\right]$$

Moreover, we shall work under the following "Predictable Representation Property" and assume that our utility functionals satisfy a Markov condition.

**Assumption 4.1.** The *Predictable Representation Property (PRP)* holds: for some $D \in \mathbb{N} \cup \{0\}$ there exists processes $w^{d+1}, ..., w^D$ adapted to the filtration $\{\mathcal{F}_t\}$ generated by the process $\bar{w}$ such that the extended process $w = (\bar{w}^1, ..., \bar{w}^d, w^{d+1}, ..., w^D)$ has uncorrelated increments which are independent from the past, have zero mean, non-trivial finite second moments, and

$$(19) \qquad L^0(\mathcal{F}_{t+1}) = \left\{ x + Z\Delta w_{t+1} : x \in L^0(\mathcal{F}_t), Z \in [L^0(\mathcal{F}_t)]^D \right\}.$$

We stress that if initially the $d$-dimensional $\bar{w}$ process driving the price had not enjoyed the PRP, then Assumption 4.1 simply says that we can complete the former process in such a way that the enlarged process does enjoy the PRP, without changing the informational structure of the model. The following example clarifies our PRP assumption.

**Example 4.2** (Bernoulli Walk). Consider in $\mathbb{R}^d$, $d$ independent Bernoulli walks $w^1, \ldots, w^d$ on the time grid $\{0, h, 2h, \ldots, T\}$ starting at 0, such that $\mathbb{P}(\Delta w_t^i = \sqrt{h}) = \mathbb{P}(\Delta w_t^i = -\sqrt{h}) = \frac{1}{2}$. They do not necessarily fulfilll (19), unless $d = 1$. Yet it is well-known that for $D = 2^d - 1$, there exists an adapted family $w^{d+1}, \cdots, w^D$ of likewise distributed random walks, such that the whole extended family $w^1, \ldots, w^D$ has increments uncorrelated to each other and independent from the past, and such that (19) holds.

We further restrict ourselves to preference functionals which satisfy the following Markov Property.

**Assumption 4.3.** The *generators* $g^l$ (l=a,p) defined by
$$Z \in [L^0(\mathcal{F}_t)]^D \mapsto g_t^l(Z) := U_t^l(Z\Delta w_{t+1}),$$
are *Markovian* in the sense that $g^a, g^p$ map $\mathbb{R}^D$ to $\mathbb{R}$.

If a preference functional $U$ satisfies the usual conditions and the PRP holds, then all the relevant information of $U_t$ is summarized by its generator. Clearly $g_t$ inherits from $U_t$ being null at the origin and concave. In the case that $P$ may only take a finite number of values, and by the "local property,"
$$\mathbb{1}_{Z(\cdot)=z} g_t(Z)(\cdot) = \mathbb{1}_{Z(\cdot)=z} g_t(z)(\cdot).$$

**Example 4.4.** For optimized certainty equivalents, the generator $g(x) := U_t(x\Delta w_t) = \sup_s \{s - \mathbb{E}(H(s - x\Delta w_t))\}$ clearly satisfies the markovianity assumption under the PRP.

*Remark* 4.5. Under Assumption 4.1 one could re-write the Agent's and Principal's recursions as Backward Stochastic Difference Equation in a direct way. In doing so we would replace $\Gamma$ by $\gamma\Delta w$ everywhere in Principal's problem, this having major advantages as by-product. First, one may drop the $L^0 - L^1$ upper semicontinuity assumption and simply work with the variational representations of the utility functionals. Indeed, by (26) of Proposition A.9 this would imply $L^0$ upper semicontinuity of $V$ (as in Principal's one step unconstrained problems) in the variables $(A, \gamma)$, which is all we need. As a consequence, the results of the previous section extend to e.g. every optimized certainty equivalent utility in the PRP case. We spare the reader the repetitive work of proving the above points, and instead proceed to a more explicit characterization of optimal contracts.

From the substitution $\Gamma_{t+1} - \mathbb{E}[\Gamma_{t+1}|\mathcal{F}_{t+1}] = \gamma\Delta w_{t+1}$ for some $\gamma \in [L^0(\mathcal{F}_t)]^D$ valid by the PRP assumption, we may call a tuple $(A, \beta, \gamma)$ without danger of confusion a *contract* (we shall always work with these variables under the PRP). Principal's recursion (12) and the incentive compatibility set $\mathbb{C}_t(\beta)$ may then be re-defined in terms of such tuple in an obvious way.

*Remark* 4.6. From equation (12) it becomes apparent that under the PRP and Markovianity Assumptions $h_t$ becomes a real number for all $t$. Indeed, everything in the one-step optimization problems (the $g$'s and $c$'s) is non-random when evaluated at non-random inputs, from which it suffices to consider $(A, \beta, \gamma) \in \mathbb{R}^N \times \mathbb{R} \times \mathbb{R}^D$ and maximize point-wise. This of course shows that in this case if there is an optimal contract, then the optimizer $(A, \beta, \gamma)$ is non-random.

### 4.1. **Computing optimal contract and necessary optimality conditions.**

Starting from the original formulation (12), we tackle the attainability issue without resorting immediately to the unconstrained variant. We will thus see that in fact solving this unconstrained problem is not only sufficient but necessary in a sense. Furthermore, in our present framework we will be able to write

down explicitly the optimal contract. We first derive the First Order Conditions (FOC) for Agent's and Principal's one-step problems:

**Lemma 4.7.** Assume that $g_t^p$ is once and $g_t^a, c_t$ are twice continuously differentiable, for $t \in \mathbb{T}$. Then:

(20) $$(A, \gamma) \in \mathbb{C}_t(\beta) \quad \text{if and only if} \quad \beta\mu - \nabla c_t(A) + \beta\sigma\nabla g_t^a(\gamma) = 0.$$

Moreover, given an optimal contract $\{(A_t, \beta_t, \gamma_t)\}$ for the Principal, and supposing for every time $t \in \mathbb{T}$ that the implied one-step contracts form a regular point for the corresponding constraints appearing in the r.h.s. of (20) -this is, the matrices $\left[\mu + \sigma\nabla g_t^a(\gamma_t) \mid \beta_t\sigma\nabla^2 g_t^a(\gamma_t) \mid -\nabla^2 c_t(A_t)\right] \in \mathbb{R}^{N \times (1+D+N)}$ have full range- there exists Lagrange multipliers $\lambda_t \in \mathbb{R}^N$ s.t. the following systems admit a solution:

(21) $$0 = [\beta_t\mu - \nabla c_t(A_t)] + \beta_t\sigma\nabla g_t^a(\gamma_t)$$

(22) $$0 = [\mu - \nabla c_t(A_t)] + \sigma\nabla g_t^p(\sigma' A_t - \gamma_t) - \nabla^2 c_t(A_t)\lambda_t$$

(23) $$0 = \nabla g_t^a(\gamma_t) - \nabla g_t^p(\sigma' A_t - \gamma_t) + \beta_t\nabla^2 g_t^a(\gamma_t)\sigma'\lambda_t$$

(24) $$0 = \lambda_t[\mu + \sigma\nabla g_t^a(\gamma_t)].$$

*Proof.* We omit the time index for simplicity. The identity (20) follows by differentiation and noticing that the optimization problem in $\mathbb{C}_t(\beta)$ is concave in the $A$ variable. It is also easy to see that the matrix

$$\left[\mu + \sigma\nabla g^a(\gamma) \mid \beta\sigma\nabla^2 g^a(\gamma) \mid -\nabla^2 c(A)\right] \in \mathbb{R}^{N \times (N+d+1)}$$

has as rows the gradients of the components of $\beta\mu - \nabla c_t(A) + \beta\sigma\nabla g_t^a(\gamma)$. By e.g. [4, Chapter 3] we thus have the existence of a Lagrange multiplier $\lambda$. Forming the Lagrangian

$$L = [A\mu - c_t(A)] + g_t^a(\gamma) + g_t^p(A \cdot \sigma - \gamma) + \lambda \cdot \{[\beta\mu - \nabla c_t(A)] + \beta\sigma\nabla g^a(\gamma)\}$$

and taking the partial derivatives w.r.t. $\lambda, A, \gamma, \beta$ yields the desired system. □

Dropping the time index again, notice that multiplying (21) by $\lambda$ yields $\lambda\nabla c(A) = 0$. Thus multiplying (23) by $\lambda'\sigma$, (22) by $\lambda$, adding them up and then multiplying by $\beta$ yields:

$$\beta\lambda'[\beta\sigma\nabla^2 g^a(\gamma)\sigma' - \nabla^2 c(A)]\lambda = 0.$$

Therefore, as soon as one searches for a $\beta > 0$ and either $c$ or $g^a$ are respectively strictly convex or concave, then necessarily $\lambda = 0$. This shows that a reasonable optimal solution to the problem must necessarily solve also the "unconstrained" problem with FOC:

$$0 = [\beta\mu - \nabla c_t(A)] + \beta\sigma\nabla g^a(\gamma)$$
$$0 = [\mu - \nabla c_t(A)] + \sigma\nabla g^p(\sigma' A - \gamma)$$
$$0 = \nabla g^a(\gamma) - \nabla g^p(\sigma' A - \gamma).$$

We knew from previous sections, in greater generality, that solving the unconstrained problem is sufficient to construct a solution to the original constrained one. Hence these last equations show that, in the present context at least, passing through the unconstrained formulation is actually also necessary, at least for contracts with $\beta > 0$.

Subtracting the second from the first equation above and then using the third one, we get:

$$(\beta - 1)[\mu + \sigma\nabla g^a(\gamma)] = 0.$$

Thus either $\beta = 1$ is optimal, or $\mu + \sigma\nabla g^a(\gamma) = 0$. This last case can be called degenerate, since under it we derive from (21) that it is optimal for the Agent to exercise minimum effort: $\nabla c(A) = 0$. Since necessary conditions give a larger set of potential optimal points than the actual set of optima, we are inclined to say that this degenerate case is suboptimal.

4.2. **Base preferences.** We close this section with an analysis of the benchmark case where both parties' preferences originate from a common base preference functional: $U^l(\cdot) = \frac{1}{\gamma^l}U(\gamma^l\cdot)$ for $l = a, p$. In terms of generators, this means that $g^l(\cdot) = \frac{1}{\gamma^l}g(\gamma^l\cdot)$. We assume that $\nabla g$ is injective. Then $(A^*, \gamma^*, \beta^*)$ satisfies the system in Lemma 4.7, with $\lambda = 0$, where $A^*$ solves

$$(25) \qquad 0 = [\mu - \nabla c_t(A^*)] + \sigma\nabla g\left[\frac{\gamma^a\gamma^p}{\gamma^a + \gamma^p}\sigma'A^*\right],$$

and $\gamma^* = \frac{\gamma^p}{\gamma^a+\gamma^p}\sigma'A^*$ and $\beta^* = 1$.

The final part of the following proposition shows to what extend the structure of optimal contracts in [20, Theorem 1] can be recovered.

**Proposition 4.8.** Under the Markovianity and PRP Assumptions, the optimal contract (interpreted as a mapping between strategies to payments) is of the form of:

$$A \mapsto \bar{S}(A) = \kappa + \sum \gamma_t^*\Delta w_{t+1} + [W_T^A - \tilde{W}_T],$$

where $W_T^A = W_0 + \sum A_t\Delta\tilde{P}_{t+1}$, $\tilde{W} = W^{A^*}$, and $\kappa \in \mathbb{R}$. Here $A^*$ and $\gamma^*$ (both vector/scalar valued deterministic processes) are the optimal ones for the Principal. Moreover, if the utilities stem from a common base functional, then we can write the optimal contract as:

$$A \mapsto \bar{S}(A) = \bar{\kappa} + \frac{\gamma^p}{\gamma^p + \gamma^a}W_T^A + \frac{\gamma^a}{\gamma^p + \gamma^a}[W_T^A - \tilde{W}_T],$$

having the form of cash plus a convex combination of the wealth generated by the Agent and the performance (gains/losses) obtained w.r.t. a benchmark portfolio, as in [20, Theorem 1].

*Proof.* By Theorem 2.15 we get:

$$\Theta = R + \sum\left[\gamma_t^*\Delta w_{t+1} + c_t(A_t^*) - A_t^*\Delta\tilde{P}_{t+1} - g_t^a(\gamma_t^*)\right] = \kappa + \sum\gamma_t^*\Delta w_{t+1} - \tilde{W}_T,$$

where we used that $\gamma^*$ and $A^*$ are optimal (Lemma 4.7), that $\beta = 1$ is optimal, that $\kappa := R + \sum c(A_t^*) - g_t^a(Z_t^a + \sigma'A_t^*)$ is a constant, thanks to Assumption 4.3, and the fact that $A_t^*$ and $\gamma_t^*$ are deterministic (Remark 4.6). Again by Theorem 2.15 this shows that the contract $A \mapsto \kappa + \sum\gamma_t^*\Delta w + W_T^A - \tilde{W}_T$ is optimal. If further the utility functionals are a re-scaling of one another, we know that $\gamma_t^* = \frac{\gamma^p}{\gamma^p+\gamma^a}\sigma'A_t^*$. Plugging in this into the previous expression for the optimal contract, we conclude. $\square$

**Example 4.9** (1d-Bernoulli Setting, Entropic Utility). Suppose Agent's and Principal's utility functions are respectively

$$U_t^a(X) = -\frac{1}{\gamma^a}\log\left(\mathbb{E}\left[e^{-\gamma^a X}|\mathcal{F}_t\right]\right) \text{ and } U_t^p(X) = -\frac{1}{\gamma^p}\log\left(\mathbb{E}\left[e^{-\gamma^p X}|\mathcal{F}_t\right]\right),$$

with $\gamma^a, \gamma^p > 0$, and that Agent's cost function is $c(a) = h\frac{a^2}{2}$. Assume also a one dimensional market driven by a simple Bernoulli-walk setting (that is $N = d = 1$: one asset, one source of randomness); see example 4.2. We first observe that $g_t(x) = -\log\left(\frac{e^{\sqrt{h}x}+e^{-\sqrt{h}x}}{2}\right) = -log \circ \cosh(\sqrt{h}x)$, from which $\nabla g_t(x) = -\sqrt{h}\tanh(\sqrt{h}x)$. From here, and manipulating (25), we get that the optimal action $A_t^*$ at time $t$ is the solution to the equation:

$$-\frac{\gamma^a\gamma^p}{\gamma^a + \gamma^p}\sqrt{h}\sigma A_t^* = \frac{1}{2}\log\left(\frac{\sigma\sqrt{h} + A_t^*h - \mu}{\sigma\sqrt{h} - A_t^*h + \mu}\right).$$

## 5. Conclusion

The present article clarifies the structure of optimal linear contracts in dynamic models of portfolio delegation when both parties' preferences satisfy translation invariance, time consistency and certain regularity conditions. We have shown how the problem of dynamic contracting can be reduced to a recursive sequence of one-period conditional optimization problems. Using conditional analysis techniques we established general attainability results for the Agent and Principal problems and derived the representation of optimal contracts found in [20] under a Markov-PRP assumption and for base preferences and general costs. Several questions are still open. First, the restriction to linear contracts is undesirable. Unfortunately, our method does in no obvious way carry over to non-linear contracts. Second, in the PRP framework we assumed that the Principal observes the driving process $\bar{w}$. Although this assumption seems common in the literature, it would be more natural to assume that the Principal observes the price increments only. This would add an additional adverse selection component to our model, if one interprets the Agent's additional information as his type, and hence leading to very complex optimization problems. Finally, it would be interesting to analyze portfolio delegation models under limited liability. If one restricts oneself a-priori to a particular class of pay-off profiles such as call options, then our methods can probably still be used to establish existence of optimal contracts (within the pre-specified class). It is an open questions how to analyze models of limited liability without any such a-priori restriction.

## Appendix A. Conditional Analysis

This appendix recalls conditional analysis results needed to analyze our dynamic contracting problem. We also establish new results which are key to our PA problem. For a detailed discussion of finite dimensional conditional analysis we refer to [7] and references therein; for a thorough treatment of conditional analysis on $L^p$ spaces we refer to [15].

A.1. **Finite dimensional conditional analysis.** On a given probability space $(\Omega, \mathcal{F}, \mathbb{P})$ we denote by $L$ and $L^0$ the sets of all, respectively all a.s. finite random variables. We apply almost-sure identification and ordering on this sets and put $\overline{L} := \{X \in L : X > -\infty\}$ and $\underline{L} := \{X \in L : X < \infty\}$ and denote by $\mathbb{N}(\mathcal{F})$ the set of variables in $L^0$ which take values in $\mathbb{N}$. We fix $N \in \mathbb{N}$ and view $E := [L^0(\mathcal{F})]^N$ as a finite-dimensional topological $L^0(\mathcal{F})$-module over the ring $L^0(\mathcal{F})$. On $E$ we define the *conditional norm* $\|X\| = (XX)^{\frac{1}{2}}$ (notice that this is a random variable), where the product is the euclidean one.

**Definition A.1.** A set $C \subset E$ is called:

- stable if $\mathbb{1}_A X + \mathbb{1}_{A^c} Y \in C$, for every $X, Y \in C$, $A \in \mathcal{F}$
- $\sigma-$stable if $\sum_{n \in \mathbb{N}} \mathbb{1}_{A_n} X_n \in C$, for every sequence $(X_n) \subset C$ and partition $(A_n) \subset \mathcal{F}$ of $\Omega$
- $L^0-$convex if $\lambda X + (1 - \lambda)Y \in C$, for every $X, Y \in C$ and $\lambda \in L^0$ with values in $[0, 1]$
- sequentially closed if it contains all the limits of its a.s. converging sequences.
- $L^0-$bounded if $\operatorname{ess\,sup}_{X \in C} \|X\| \in L^0$.

A stable and sequentially closed set is $\sigma-$stable. We define for $M \in \mathbb{N}(\mathcal{F})$ and $(X_n) \subset E$ the element $X_M = \sum_{n \in \mathbb{N}} \mathbb{1}_{M=n} X_n \in E$ and notice that if the former sequence belongs to a $\sigma-$stable set, then the latter does so too. The following result is a generalization of the classical Bolzano-Weierstrass Theorem.

**Lemma A.2.** Let $(X_n) \subset E$ be $L^0-$bounded. Then there exists $X \in E$ and a sequence $(N_n) \in \mathbb{N}(\mathcal{F})$ such that $N_{n+1} > N_n$ and $X = \lim_{n \to \infty} X_{N_n}$ a.s. Also, if $(x_n) \subset L^0$ is such that $x := \limsup x_n \in L^0$, then there exists a sequence $(N_n) \in \mathbb{N}(\mathcal{F})$ such that $N_{n+1} > N_n$ and $x = \lim_{n \to \infty} x_{N_n}$ a.s.

*Proof.* For the first statement refer to [7, Theorem 3.8]. For the second, define $N_0 = 0$ and $N_n = \min\{m > N_{n-1} : x_m \geq x - 1/n\}$. Then $N_n \in \mathbb{N}(\mathcal{F})$ and $N_{n+1} > N_n$, from which $N_n \geq n$ follows. Now, notice that $\sup_{m \geq n} x_m \geq \sup_{m \geq N_n} x_m \geq x_{N_n} \geq x - 1/n$ a.s., from which $x = \lim_{n \to \infty} x_{N_n}$ a.s. $\qquad\square$

As in the euclidean case, convexity opens the way to a necessary and sufficient characterization of boundedness (see [7, Theorem 3.13]):

**Theorem A.3.** Let $C$ be a sequentially closed $L^0$−convex subset of $E$ which contains 0. Then $C$ is $L^0$−bounded if and only if for any $X \in C\backslash\{0\}$ there exists a $k \in \mathbb{N}$ such that $kX \notin C$.

Let us now introduce the notions of continuity, convexity and stability of functions defined on subsets of $E$ and taking values in a set of random variables.

**Definition A.4.** Let $C \subset E$. A function $f : C \to L$ is called:

- $L^0$−lower semicontinuous at $X \in C$ if $f(X) \leq \liminf f(X_n)$ for every sequence $(X_n) \subset C$ with a.s. limit $X$.
- $L^0$−continuous at $X \in C$ if $f(X) = \lim f(X_n)$ whenever $(X_n) \subset C$ has a.s. limit $X$.
- $L^0$−convex if $f(\lambda X + (1 - \lambda)Y) \leq \lambda f(X) + (1 - \lambda)f(Y)$, for every $X, Y \in C$ and $\lambda \in L^0$ with values in $[0, 1]$
- stable if $f(\mathbb{1}_A X + \mathbb{1}_{A^c} Y) = \mathbb{1}_A f(X) + \mathbb{1}_{A^c} f(Y)$, for every $X, Y \in C$, $A \in \mathcal{F}$.

For the last two items it is assumed that $C$ is $L^0$−convex, respectively, stable. Strict $L^0$−convexity is defined in terms of a strict inequality. Finally $f$ is called (upper/lower semi)continuous on $C$ if it is so on every point of $C$. If $f$ is continuous and stable over a $\sigma$−stable and sequentially closed set, then it satisfies the stability property for countable partitions too. If $f$ is $L^0$−convex or $L^0$−concave, then it is local (meaning $\mathbb{1}_A f(X) = \mathbb{1}_A f(Y)$ whenever $\mathbb{1}_A X = \mathbb{1}_A Y$), which in itself directly implies that it also satisfies the stability property for countable partitions.

The following result is implied by the proof of [7, Theorem 4.13], since all the authors really use is $\sigma$−stability of the set under consideration (which is implied by their stronger assumptions). We give a self-contained proof here.

**Lemma A.5.** If a non-empty set $C \subset E$ is $\sigma$−stable and is not $L^0$−bounded, then there is a set $\tilde{\Omega}$ with $\mathbb{P}(\tilde{\Omega}) > 0$ and a sequence $\{X_n\} \subset C$ such that, for every $n \in \mathbb{N}$, $|X_n| \geq n$ over $\tilde{\Omega}$

*Proof.* We define $U_n := \{B \in \mathcal{F} : \exists X \in C, |X| \geq n \text{ on } B\}$, which is non-empty since $C$ is unbounded, introduce the family of decreasing sets $A_n := \left\{ \underset{B \in U_n}{\mathrm{ess\,sup}}\, \mathbb{1}_B = 1 \right\}$ and put $A := \bigcap_n A_n$. Assuming that $\mathbb{P}(A) = 0$, or equivalently that $\mathbb{P}(\cup_n A_n^c) = 1$, then for every $X \in C$:

$$|X| = \left| \sum_n X \mathbb{1}_{\{A_n^c \cap A_{n-1}\}} \right| \leq \sum_n |X| \mathbb{1}_{\{A_n^c \cap A_{n-1}\}} \leq \sum_n n \mathbb{1}_{\{A_n^c \cap A_{n-1}\}} \in L^0(\mathcal{F}).$$

Since $X \in C$ was arbitrary, this implies that $C$ is $L^0(\mathcal{F})$−bounded. Therefore $\mathbb{P}(A) > 0$ must hold. By definition of ess sup we have that there exist $\{B^{l,n}\}_l \in U_n$ such that $\underset{B \in U_n}{\mathrm{ess\,sup}}\, \mathbb{1}_B = \sup_l \mathbb{1}_{B^{l,n}}$ a.s. This implies $A_n = \bigcup_l B^{l,n}$ a.s. Taking $X^{l,n}$ such that $|X^{l,n}| \geq n$ on $B^{l,n}$, and fixing an $X^* \in C$ arbitrary, let us define:

$$X^{(n)} := X^* \mathbb{1}_{\left\{ (\cup_l B^{l,n})^c \right\}} + \sum_l X^{l,n} \mathbb{1}_{\left\{ B^{l,n} \cap (\cup_{m<l} B^{m,n})^c \right\}} + X^{0,n} \mathbb{1}_{B^{0,n}},$$

which belongs to $C$ thanks to $\sigma$−stability. Clearly

$$|X^{(n)}| \geq n \mathbb{1}_{\left\{ \cup_l B^{l,n} \right\}} + |X^*| \mathbb{1}_{\left\{ (\cup_l B^{l,n})^c \right\}},$$

and therefore a.s. $|\mathbb{1}_A X^{(n)}| \geq n\mathbb{1}_A$. Thus we have that $|X^{(n)}| \geq n$ on $A$ for every $n$. $\qquad\square$

The following conditional optimization theorem is used to prove attainability of the Agent Problem. For a proof we refer to [7, Theorem 4.4].

**Theorem A.6.** Let $C$ be a sequentially closed and stable subset of $E$ and $f : C \to \overline{L}$ be a $L^0-$lower semicontinuous, stable function. Assume there exists an $X_0 \in C$ such that the set $\{X \in C : f(X) \leq f(X_0)\}$ is $L^0-$bounded. Then there exists an $\hat{X} \in C$ such that

$$f\left(\hat{X}\right) = \operatorname*{ess\,inf}_{X \in C} f(X).$$

If $f$ and $C$ are $L^0-$convex then the "arg min" set is also $L^0-$convex, and if in addition $f$ is strictly $L^0-$convex then $\hat{X}$ is the sole (a.s.) optimizer.

We finally adapt a Komlos-type lemma (as in [12, Lemma A1.1]) for conditionally bounded random variables, which we use to prove our general attainability result (Theorem 2.18). We thank a referee for hinting at the proof we give now.

**Lemma A.7.** Let $\{\xi_n\}_n$ be $[0, +\infty)$-valued random variables defined on a common probability space $(\Omega, \mathcal{G}, \mathbb{P})$, take $\mathcal{F}$ a sub-sigma algebra and assume that the set $C := conv\{\xi_n : n \in \mathbb{N}\}$ satisfies the following conditional boundedness condition:

$$\forall \epsilon \in L^0_+(\mathcal{F}), \exists a \in L^0(\mathcal{F}) \text{ such that } \forall h \in C, \mathbb{P}(h \geq a|\mathcal{F}) \leq \epsilon.$$

Then there exists a $[0, +\infty)$-valued random variable $X$ and a sequence $\{x_n\}$, where $x_n$ belongs to the convex hull of $\{\xi_n, \xi_{n+1}, \dots\}$ such that $x_n \to X$ almost surely.

*Proof.* By [12, Lemma A1.1] it suffices to show that $C$ is bounded in probability. By assumption, we have that $p_n := \operatorname{ess\,sup}_{h \in C} \mathbb{P}(h \geq n|\mathcal{F}) \to 0$, as $n \to \infty$ and $\mathbb{P}-$a.s. Since also $p_n \in [0, 1]$ a.s. we conclude by dominated convergence that $\mathbb{E}[p_n] \to 0$, which of course is stronger than $\sup_{h \in C} \mathbb{P}[h \geq n] \to 0$, so we conclude. $\qquad\square$

A.2. **Conditional analysis on $L^p$.** Let $\mathcal{F}$ be a sub sigma-algebra of $\mathcal{G}$. For every $p \in [1, +\infty]$ we define:

$$||X||_p = \begin{cases} \mathbb{E}[|X|^p|\mathcal{F}] & \text{if } p \in [1, +\infty) \\ \operatorname{ess\,inf}\{Y \in L^0_+(\mathcal{F}) \text{ s.t. } Y \geq |X|\} & \text{if } p = +\infty. \end{cases}$$

This is well defined for every $X \in L^0(\mathcal{G})$. We further define the conditional $L^p$-space

$$L^p_{\mathcal{F}}(\mathcal{G}) := \{X \in L^0(\mathcal{G}) \text{ st. } ||X||_p \in L^0(\mathcal{F})\}.$$

It is shown in [15] that $L^p_{\mathcal{F}}(\mathcal{G})$ is a topological $L^0(\mathcal{F})-$module over the topological ring $L^0(\mathcal{F})$, and $|| \cdot ||_p$ is an $L^0(\mathcal{F})-$norm inducing the module topology on $L^p_{\mathcal{F}}(\mathcal{G})$.

A function $U : L^p_{\mathcal{F}}(\mathcal{G}) \to \underline{L}^0$ is called:

- $L^0(\mathcal{F})-$concave: if $U(\lambda X + (1 - \lambda)X') \geq \lambda U(X) + (1 - \lambda)U(X')$ for every $\lambda \in L^0(\mathcal{F}) \cap [0, 1]$ and every $X, X' \in L^p_{\mathcal{F}}(\mathcal{G})$
- proper: if $\exists X \in L^p_{\mathcal{F}}(\mathcal{G})$ such that $U(X) > -\infty$ and $\forall X' \in L^p_{\mathcal{F}}(\mathcal{G})$ it holds $U(X) < \infty$
- $L^p_{\mathcal{F}}(\mathcal{G})$-upper semicontinuous: if for every net $\{X_\alpha\} \subset L^p_{\mathcal{F}}(\mathcal{G})$ converging to some $X$ in conditional norm, it holds that $\operatorname{ess\,inf}_\beta \operatorname{ess\,sup}_{\alpha \geq \beta} U(X_\alpha) \leq U(X)$
- monotone: if $U(X) \geq U(X')$ whenever $X \geq X'$
- translation invariant: if $U(X + Y) = U(X) + Y$ for every $X \in L^p_{\mathcal{F}}(\mathcal{G})$ and $Y \in L^0(\mathcal{F})$

The following representation result re-phrases [15, Corollary 3.14]:

**Theorem A.8.** Let $p \in [1, \infty)$ and $U : L^p_{\mathcal{F}}(\mathcal{G}) \to \underline{L}^0(\mathcal{F})$ satisfy the above conditions. Let $q$ be the Hölder conjugate of $p$ and define

$$\mathcal{W} := \{Z \in L^q_{\mathcal{F}}(\mathcal{G}) : Z \geq 0, \mathbb{E}[Z|\mathcal{F}] = 1\}, \quad \alpha(Z) := \operatorname*{ess\,sup}_{X \in L^p_{\mathcal{F}}(\mathcal{G})} \{U(X) - \mathbb{E}[ZX|\mathcal{F}]\}.$$

Then $U$ satisfies the following variational representation:

$$U(X) = \operatorname*{ess\,inf}_{Z \in \mathcal{W}} \{\mathbb{E}[ZX|\mathcal{F}] + \alpha(Z)\}.$$

In the next Proposition we prove that $L^p_{\mathcal{F}}(\mathcal{G})$−upper semicontinuity is a consequence of $L^0 - L^p$ upper semicontinuity (see Definition 2.6). This of course implies Proposition 2.8.

**Proposition A.9.** Let $U : L^p_{\mathcal{F}}(\mathcal{G}) \to \underline{L}^0(\mathcal{F})$ be $L^0 - L^p$ upper semicontinuous. Then $U$ is also $L^p_{\mathcal{F}}(\mathcal{G})$−upper semicontinuous. Furthermore, if $U$ is also proper, monotone, translation invariant and $L^0(\mathcal{F})$-concave, then $U$ admits a variational representation and for any $N \in \mathbb{N}$ and $\Delta \in [L^p_{\mathcal{F}}(\mathcal{G})]^N$ the functional

(26)
$$A \in [L^0(\mathcal{F})]^N \mapsto U(A\Delta)$$

is $L^0$-upper semicontinuous in the sense of Definition A.4. Under the same hypotheses, if $A_n \in [L^0(\mathcal{F})]^N \to A$ a.s. and $\{\Gamma^n\}_n$ is $L^p_{\mathcal{F}}(\mathcal{G})$-bounded such that $\Gamma^n \to \Gamma$ a.s. then

$$\limsup_n U(A_n\Delta + \Gamma_n) \leq U(A\Delta + \Gamma)$$

*Proof.* For the first part, by [14, Lemma 3.10], it is enough to prove that the sets $K_k := \{X \in L^p_{\mathcal{F}}(\mathcal{G}) : U(X) \geq k\}$ are conditionally closed for every $k \in L^0(\mathcal{F})$. We will prove that their complements are conditionally open. To this end we fix such a $k$ and and assume to the contrary that $(K_k)^c$ is not open. We thus take $X$ such that $U(X) < k$ on a non-negligible set and such that for every $N \in \mathbb{N}(\mathcal{F})$ we have that $K_k \cap B(X, 1/N) \neq \emptyset$, where $B(X, 1/N) = \{Z : \mathbb{E}(|Z - X|^p|\mathcal{F}) \leq 1/N\}$. This means that we can find, for every $N \in \mathbb{N}(\mathcal{F})$, an element $X_N \in B(X, 1/N)$ such that $U(X_N) \geq k$ a.s. A straightforward adaptation of Markov's inequality yields

$$\mathbb{P}(|X_N - X| \geq \epsilon|\mathcal{F}) \leq \frac{\mathbb{E}(|X_N - X|^p|\mathcal{F})}{\epsilon^p}$$

for every $\epsilon \in L^0(\mathcal{F})_{++}$. From this we may find for every natural number $n$ an element $M_n \in \mathbb{N}(\mathcal{F})$ such that:

- for every $N \in \mathbb{N}(\mathcal{F})$ st. $N \geq M_n$ it holds that $\mathbb{P}(|X_N - X| \geq 1/n|\mathcal{F}) \leq 1/n^2$ a.s.
- for every n: $M_{n+1} > M_n$ a.s.

Now, we will use a "Borel-Cantelli Lemma"-type reasoning in order to prove that the sequence $\{X_{M_n}\}$ converges almost surely to $X$. First notice that for a fixed $l \in \mathbb{N}$:

$$\sum_{n \in \mathbb{N}} \mathbb{P}(|X_{M_n} - X| \geq 1/l|\mathcal{F}) \leq \sum_{n \leq l} \mathbb{P}(|X_{M_n} - X| \geq 1/l|\mathcal{F}) + \sum_{n > l} \mathbb{P}(|X_{M_n} - X| \geq 1/n|\mathcal{F}),$$

and since the last term is bounded above by $\sum_{n>l} 1/n^2$, the original sum belongs to $L^0(\mathcal{F})$ (and so is a.s. finite). Define now $i.o. \{|X_{M_.} - X| \geq 1/l\} := \bigcap_{m \in \mathbb{N}} \bigcup_{n \geq m} \{|X_{M_n} - X| \geq 1/l\}$. Then:

$$\mathbb{P}\left(i.o. \{|X_{M_.} - X| \geq 1/l\}|\mathcal{F}\right) \leq \mathbb{P}\left(\bigcup_{n \geq m} \{|X_{M_n} - X| \geq 1/l\}|\mathcal{F}\right) \leq \sum_{n \geq m} \mathbb{P}(|X_{M_n} - X| \geq 1/l|\mathcal{F}),$$

and so the left-hand side does not depend on $m$ whereas the right one tends a.s. to 0 as $m$ increases. This shows that $\mathbb{P}(i.o. \{|X_{M_.} - X| \geq 1/l\}|\mathcal{F}) = 0$ a.s. Taking expectations, $\mathbb{P}(i.o. \{|X_{M_.} - X| \geq 1/l\}) = 0$. Since this holds for every $l$, we conclude that indeed $\{X_{M_n}\}$ converges almost surely to $X$.

Finally we have by the $L^0 - L^p$ upper semicontinuity assumption that $k \leq \limsup_n U(X_{M_n}) \leq U(X)$ a.s. since by definition $X_{M_n} \in B(X, 1)$, but also $U(X) < k$ on a non-negligible set, which is a contradiction. This completes the proof of the first statement.

By Theorem A.8 and the first claim we know that $U$ has a variational representation. That $A \mapsto U(A\Delta)$ is $L^0$-upper semicontinuous is a consequence of the last claim in the proposition (taking $\Gamma_n = 0$). So to establish the last claim and finish the proof, is suffices to compute

$$\mathbb{E}\left[|A_n\Delta + \Gamma_n|^p|\mathcal{F}\right]^{1/p} \leq C \sup_{i=1,\ldots,N} |A_n^i|\mathbb{E}\left[|\Delta|^p|\mathcal{F}\right]^{1/p} + \mathbb{E}\left[|\Gamma_n|^p|\mathcal{F}\right]^{1/p},$$

and observe that the r.h.s. is bounded from above by some r.v. in $L^0(\mathcal{F})$, by conditional $L^p$-boundedness of the $\Gamma_n$ and since the (components of) the $A_n$ converge a.s. All in all $\{A_n\Delta + \Gamma_n\}_n$ is $L^p_{\mathcal{F}}(\mathcal{G})$-bounded and converges a.s. to $A\Delta + \Gamma$, so we conclude by the $L^0 - L^p$ upper semicontinuity assumption. $\square$

## APPENDIX B. OPTIMIZED CERTAINTY EQUIVALENTS AND PROOF OF THEOREM 2.21

We start with a number of technical results for the examples in Section 2.2.

**Lemma B.1.** The following hold.

(i) The extensions in (3) are well defined for TVAR and, more generally, for optimized certainty equivalent families for which both $1 \in \cap_t int(dom(H_t^*))$ and every $H_t$ is bounded from below (equivalently $0 \in dom(H_t^*)$).

(ii) Take $\mathcal{F} = \mathcal{F}_t$, $\mathcal{G} = \mathcal{F}_{t+1}$ for $t$ fixed. For any $\gamma > 0$ and $\lambda \in (0, 1)$, the entropic functional

$$X \in L^1_{\mathcal{F}}(\mathcal{G}) \mapsto -\frac{1}{\gamma} \log\left(\mathbb{E}(exp(-\gamma X)|\mathcal{F})\right),$$

as well as the Tail-value-at-risk functional

$$X \in L^1_{\mathcal{F}}(\mathcal{G}) \mapsto \operatorname*{ess\,sup}_s \left\{s - \lambda^{-1}\mathbb{E}([s - X]_+|\mathcal{F})\right\},$$

are $L^0 - L^1$ u.s.c. More generally, optimized certainty equivalents for which $1 \in int(dom(H^*))$ are $L^0 - L^1$ u.s.c.

(iii) The TVAR family and, more generally, OCE families for which $1 \in \cap_t int(dom(H_t^*))$ and every $H_t$ is bounded from below, all satisfy Assumption 2.7 after pasting.

*Proof.* For TVAR we have $H(l) = \lambda^{-1}[l]_+$ and $H^*(x) = \Psi_{[0,\lambda^{-1}]}(x)$, the convex indicator of $[0, \lambda^{-1}]$. In particular, $1 \in int(dom(H^*))$ and $H$ is lower bounded. In the following we work with abstract OCEs for which the latter conditions hold.

i) It suffices to show that the extensions defined by (3) produce a functional which never attains the value $+\infty$. Let $U$ stand for any OCE (associated to $H$) satisfying the stated properties and let $X$ be a r.v. not attaining $+\infty$. For $1 + \epsilon \in int(dom(H^*))$ we define

$$N^\epsilon := \sum_{n=1}^\infty n\mathbb{1}_{\mathbb{P}(X \leq n|\mathcal{F}) > [1+\epsilon]^{-1}, \mathbb{P}(X \leq n-1|\mathcal{F}) \leq [1+\epsilon]^{-1}} = \inf\{n \in \mathbb{N} : \mathbb{P}(X \leq n|\mathcal{F}) > [1+\epsilon]^{-1}\},$$

which is then finite and belongs to $\mathbb{N}(\mathcal{F})$. Inspired by [8, Proof of (12)] we introduce a partition $A_0 := \{\mathbb{P}(X \leq N^\epsilon|\mathcal{F}) > 0\}$ and $A_n := \{\mathbb{P}(X \leq N^\epsilon + n|\mathcal{F}) > 0, \mathbb{P}(X \leq N^\epsilon + n - 1|\mathcal{F}) = 0\}$ for $n \geq 1$, so we define

$$\xi = \sum_{n \geq 0} \mathbb{1}_{A_n} \frac{\mathbb{1}_{X \leq N^\epsilon + n}}{\mathbb{P}(X \leq N^\epsilon + n|\mathcal{F})}.$$

It is then easy to see that $\mathbb{E}[\xi|\mathcal{F}] = 1$ and that $\xi \in [0, 1+\epsilon]$, so that a.s. $\xi \in dom(H^*)$. We conclude that

$$
\begin{aligned}
U(X \wedge k) &= \underset{s}{\text{ess sup}}\{s - \mathbb{E}[H(s - X \wedge k)]\} \\
&\leq \underset{s}{\text{ess sup}}\{s - s\mathbb{E}[\xi|\mathcal{F}] + \mathbb{E}[\xi(X \wedge k) + H^*(\xi)|\mathcal{F}]\} \\
&\leq \mathbb{E}[\xi X|\mathcal{F}] + \mathbb{E}[H^*(\xi)|\mathcal{F}] \leq \sum_{n \geq 0} \mathbb{1}_{A_n}[N^\epsilon + n] + \mathbb{E}[H^*(\xi)|\mathcal{F}] < +\infty,
\end{aligned}
$$

where finiteness comes from the fact that $H^*$ must send $[0, 1 + \epsilon]$ into a bounded set. Hence $\lim_{k \to +\infty} U(X \wedge k) < +\infty$ and we conclude.

ii) Let us take $X_n$ bounded in $L^1_{\mathcal{F}}(\mathcal{G})$ such that $X_n \to X$ a.s. For any $C \in L^0(\mathcal{F})$ we want to show that if $\text{ess sup}_s \{s - \mathbb{E}(H(s - X_n)|\mathcal{F})\} \geq C$ then also $\text{ess sup}_s \{s - \mathbb{E}(H(s - X)|\mathcal{F})\} \geq C$. Indeed, let us first take $s_n \in L^0(\mathcal{F})$ such that

$$
s_n - \mathbb{E}(H(s_n - X_n)|\mathcal{F}) \geq C - n^{-1}.
$$

Because $H$ is convex, lower-semicontinuous and proper, we have that $H(s_n - \mathbb{E}(X_n|\mathcal{F})) \geq R[s_n - \mathbb{E}(X_n|\mathcal{F})] - H^*(R)$ for each $R$, and so in particular

$$
s_n[1 - R] \geq C - n^{-1} - H^*(R) - R\mathbb{E}(X_n|\mathcal{F}),
$$

for every $R \in dom(H^*)$. If we were able to find $R \in dom(H^*) \cap (1, \infty)$ then for such element we would have

$$
s_n \leq \frac{C - n^{-1} - H^*(R) - R\mathbb{E}(X_n|\mathcal{F})}{1 - R}.
$$

Similarly, if $r \in dom(H^*) \cap (-\infty, 1)$ existed then we would get

$$
s_n \geq \frac{C - n^{-1} - H^*(r) - r\mathbb{E}(X_n|\mathcal{F})}{1 - r}.
$$

Altogether, we could conclude that the quantities $s_n$ are $L^0(\mathcal{F})$-bounded, since the random variables $X_n$ were bounded in $L^1_F(\mathcal{G})$. By Lemma A.2 we could find $N_n \in \mathbb{N}(\mathcal{F})$ increasing to $+\infty$ such that $s_{N_n} \to \bar{s}$ a.s. for some $\bar{s} \in L^0(\mathcal{F})$, and obviously $X_{N_n} \to X$ a.s. still. By locality we would have that

$$
s_{N_n} - \mathbb{E}(H(s_{N_n} - X_{N_n})|\mathcal{F}) \geq C - N_n^{-1}.
$$

Taking $\limsup_n$, using the fact that $H$ is bounded below by an affine function and conditional Fatou's Lemma, we may obtain

$$
\bar{s} - \mathbb{E}(H(\bar{s} - X)|\mathcal{F}) \geq C.
$$

This readily implies what we wanted to prove. To conclude, observe that the conditions $dom(H^*) \cap (1, \infty) \neq \emptyset$ and $dom(H^*) \cap (-\infty, 1) \neq \emptyset$ are together equivalent to $1 \in int(dom(H^*))$ in our case, since $1 \in dom(H^*)$ by definition and $H^*$ is convex.

iii) Upper semincontinuity has already been dealt with. Under the stated conditions pasting is also already justified, so it only remains to show the condition on the domain of $U_t$. If $X \in dom(U_t)$, then there must be some $s \in \mathbb{R}$ such that $\mathbb{E}[H(s - X)|\mathcal{F}_t] < \infty$. But as $H$ is bounded from below, this implies $\mathbb{E}[H(s + X_-)|\mathcal{F}_t] < \infty$, and by the normalization property on $H$ (implying $H(l) \geq l$) we get in turn $\mathbb{E}[X_-|\mathcal{F}_t] < \infty$ as desired.

$\square$

*Remark* B.2. The reason we had to prove [8, Proof of (12)] anew in the first part of the preceding proof, is that we want to include the case where $dom(H^*) \neq [0, \infty)$, in order to cover e.g. the TVAR family. This creates the difficulty of finding a $\xi$ satisfying simultaneously $\mathbb{E}[\xi|\mathcal{F}] = 1$, $\xi \in [0, 1 + \epsilon]$ and $\mathbb{E}[\xi X|\mathcal{F}] < \infty$.

*Remark* B.3. In case $U_t^a$ arises as an optimized certainty equivalent (see example 2.4), we know by Lemma B.1 and Proposition A.9 that it has a variational representation. Notice then by Young inequality that

$$\alpha_t(Z) \leq \operatorname*{ess\,sup}_{X} \operatorname*{ess\,sup}_{s} \{ s - \mathbb{E}[H_t(s - X)|\mathcal{F}_t] - \mathbb{E}[ZX|\mathcal{F}] \} \leq \mathbb{E}[H_t^*(Z)|\mathcal{F}_t],$$

whenever $Z \in L^\infty_{\mathcal{F}_t}(\mathcal{F}_{t+1})$ is s.t. $\mathbb{E}[Z|\mathcal{F}_t] = 1$. In [8] and [9] it is proved that under given conditions there is equality above (for $Z \in \mathcal{W}_t$). In any case we see that the conditions in Theorem 2.18 on $K^{a,p}$ are satisfied if these utility functionals are such that $1 \in int(dom(H_t^*))$; indeed, we may just take $\epsilon > 0$ such that $[1 - \epsilon, 1 + \epsilon] \subset dom(H_t^*)$.

Equipped with the previous result, we provide the proof of our general existence of optimal contracts.

*Proof of Theorem 2.21.* We may assume $\mathcal{F} = \mathcal{F}^A$. Under the given condition on the $H$'s, the condition on the $K$'s (see Theorem 2.18) is satisfied, thanks to Remark B.3. By Proposition 2.20 it remains to show that $h_{t+1} \in dom(U_t^p)$ for each $t$.

Either using Young's inequality or invoking Remark B.3, we know that $U_t^a(X), U_t^p(X) \leq \mathbb{E}[X|\mathcal{F}_t]$. From this we see

$$V_t(A, \Gamma) \leq -c(A) + A\mathbb{E}[\Delta\tilde{P}|\mathcal{F}_t] + \mathbb{E}[h_{t+1}|\mathcal{F}_t]$$

$$\leq -\left( \frac{c_t(A)}{|A|} - \frac{2p_+}{p_-} \right)|A| + \mathbb{E}[h_{t+1}|\mathcal{F}_t]$$

$$\leq K_t + \mathbb{E}[h_{t+1}|\mathcal{F}_t],$$

where we used that $|\Delta\tilde{P}_{t+1}| = |diag(P_t)^{-1}\Delta P_{t+1}| \leq \frac{2p_+}{p_-}$ and the existence of $K_t \in \mathbb{R}$ is a consequence of the growth of $c_t$ and its continuity. From here we get by definition that $h_t \leq K_t + \mathbb{E}[h_{t+1}|\mathcal{F}_t]$ and since $h_T = 0$ by backwards inductions follows that $h_{t+1}^p \leq L$ for some constant $L$ and all $t$. Monotonicity of the $U^p$'s mean that $h_{t+1} \in dom(U_t^p)$. $\square$

## APPENDIX C. PROOF OF THEOREM 2.15

First we turn our attention to the Agent's recursion. Let $\bar{a}$ be a generic sequence of efforts. From equation (7), we see that defining

$$H_t(\bar{a}_t, \ldots, \bar{a}_{T-1}) := U_t^a\left( \Theta(P_{0:T}) + \sum_{s \geq t} \{\beta_s \Delta W_{s+1}^{\bar{a}} - c_{s+1}(\bar{a}_{s+1})\} \right) - c_t(\bar{a}_t),$$

we get the recursion

$$H_T(\bar{a}_t, \ldots, \bar{a}_{T-1}) = \Theta(P_{0:T}),$$

$$H_t(\bar{a}_t, \ldots, \bar{a}_{T-1}) = U_t^a\left( H_{t+1}(\bar{a}_{t+1}, \ldots, \bar{a}_{T-1}) + \beta_t \bar{a}_t \Delta\tilde{P}_{t+1} \right) - c_t(\bar{a}_t).$$

Then, in terms of $H_t := \operatorname*{ess\,sup}_{a_t, \ldots, a_{T-1}} H_t(a_t, \ldots, a_{T-1})$, we get:

$$H_t(\bar{a}_t, \ldots, \bar{a}_{T-1}) \leq -c(\bar{a}_t) + U_t^a\left( \operatorname*{ess\,sup}_{a_{t+1}, \ldots, a_{T-1}} H_{t+1}(a_{t+1}, \ldots, a_{T-1}) + \beta_t \bar{a}_t \Delta\tilde{P}_{t+1} \right).$$

This yields that $H_t \leq \operatorname*{ess\,sup}_{a_t} \{ -c(a_t) + U_t^a\left( H_{t+1} + \beta_t a_t \Delta\tilde{P}_{t+1} \right) \}$. For $t = T-1$ this is an equality and by assumption the value $H_{T-1}$ is attained at some $\hat{a}_{T-1}$. Suppose now that equality holds in the previous

equation for $t+1, \ldots, T-1$, and $H_{t+1}$ was attained say at $(\hat{a}_{t+1}, \ldots, \hat{a}_{T-1})$. This implies that:

$$
\begin{aligned}
H_t &\leq \operatorname*{ess\,sup}_{a_t} \left\{ -c(a_t) + U_t^a \left( H_{t+1}(\hat{a}_{t+1}, \ldots, \hat{a}_{T-1}) + \beta_t a_t \Delta \tilde{P}_{t+1} \right) \right\} \\
&\leq \operatorname*{ess\,sup}_{a_t, \ldots, a_{T-1}} \left\{ -c(a_t) + U_t^a \left( H_{t+1}(a_{t+1}, \ldots, a_{T-1}) + \beta_t a_t \Delta \tilde{P}_{t+1} \right) \right\} \\
&= \operatorname*{ess\,sup}_{A_t, \ldots, A_{T-1}} H_t(A_t, \ldots, A_{T-1}) =: H_t.
\end{aligned}
$$

So indeed at time $t$ also $H_t = \operatorname*{ess\,sup}_{a_t} \left\{ -c(a_t) + U_t^a \left( H_{t+1} + \beta_t a_t \Delta \tilde{P}_{t+1} \right) \right\}$ holds and by assumption the last term is attained at some $\hat{a}_t$, from which $H_t$ is attained at $(\hat{a}_t, \ldots, \hat{a}_{T-1})$. This closes the inductive step, and therefore the desired recursion holds.

Now we will establish rigorously recursion (11) (equivalently (9)). To this end we denote by $\beta = (\beta_t)_t$ a generic decision variable for the Principal and $a = (a_t)_t$ where $a_t \in L^0(\mathcal{F}_t)^N$, a corresponding optimal effort for the Agent. Let

$$
N := \sum_{s \geq t+1} \left[ (1 - \beta_s) a_s \Delta \tilde{P}_{s+1} - \Delta H_{s+1} \right].
$$

Then using the just proven expression for $H_t$ (i.e. (8)), and setting $\Gamma = \beta_t a_t \Delta \tilde{P}_{t+1} + H_{t+1}$, we get:

$$
\begin{aligned}
U_t^p \left( \sum_{s \geq t} \left[ (1 - \beta_s) a_s \Delta \tilde{P}_{s+1} - \Delta H_{s+1} \right] \right) &= U_t^p \left( (1 - \beta_t) a_t \Delta \tilde{P}_{t+1} - H_{t+1} - c_t(a_t) + N \right. \\
&\qquad \left. + U_t^a \left( H_{t+1} + \beta_t a_t \Delta \tilde{P}_{t+1} \right) \right), \\
&= U_t^p \left( a_t \Delta \tilde{P}_{t+1} - \Gamma + U_t^a(\Gamma) - c_t(a_t) + N \right), \\
&= U_t^a(\Gamma) - c_t(a_t) + U_t^p \left( a_t \Delta \tilde{P}_{t+1} - \Gamma + N \right).
\end{aligned}
$$

And now, applying time-consistency and translation invariance in the last term above we get:

$$
U_t^p \left( \sum_{s \geq t} \left[ (1 - \beta_s) a_s \Delta \tilde{P}_{s+1} - \Delta H_{s+1} \right] \right) = U_t^a(\Gamma) - c_t(a_t) + U_t^p \left( a_t \Delta \tilde{P}_{t+1} - \Gamma + U_{t+1}^p(N) \right).
$$

Therefore calling $h_{t+1}(a, \Gamma) = U_{t+1}^p(N)$, we obtain recursion (11). That is to say, if $(a, \Gamma)$ does not satisfy this recursion, they will not be chosen by the Principal. In the same way we conclude for $(a, \beta)$ and recursion (9). With these recursions for $h_t(\cdot)$ already established, we can proceed to prove (12) the same way we proved the recursion for $H$. First recall that actually $h_t(a, \Gamma)$ is a short-hand for $h_t((a_s, \Gamma_{s+1})_{s \geq t})$. From this and (11) we have:

$$
h_t((\bar{a}_s, \bar{\Gamma}_{s+1})_{s \geq t}) \leq U_t^p \left( \operatorname*{ess\,sup}_{A, \Gamma} h_{t+1}(A, \Gamma) + \bar{a}_t \Delta \tilde{P}_{t+1} - \bar{\Gamma}_{t+1} \right) + U_t^a(\bar{\Gamma}_{t+1}) - c_t(\bar{a}_t).
$$

This yields

$$
h_t \leq \operatorname*{ess\,sup}_{(a_t, \Gamma_{t+1}) \in \mathbb{C}_t(\beta_t)} U_t^p \left( h_{t+1} + a_t \Delta \tilde{P}_{t+1} - \Gamma_{t+1} \right) + U_t^a(\Gamma_{t+1}) - c_t(a_t).
$$

For $t = T - 1$ this is an equality (we defined $h_T = 0$) and by assumption the value $h_{T-1}$ is attained. Using induction, similarly as how we did it for $H$, we get that (12) holds.

The validity of the change of variables $\beta_t a_t \Delta \tilde{P}_{t+1} + H_{t+1} \to \Gamma_{t+1}$ and the introduction of $\mathbb{C}(\beta)$ as a constraint inducing incentive compatibility are now obvious. This means that $h$ represents the future wealth prospects of the Principal. Hence at time $t = 0$ we obtain a solution for the whole Principal's problem, proving as well that Principal's optimal wealth is $W_0 - R + h_0$.

We proceed now to prove that a solution to Principal's recursion delivers indeed an optimal (dynamic) contract, and that the Agent behaves as predicted. Call $(\beta_t, A_t, \Gamma_{t+1})_t$ the optimal quantities attaining $h$ in (12). Define $\Theta$ and the contract $S$ as in the statement of the present theorem. Then:

$$
\begin{aligned}
U_{T-1}^a\left(\Theta + \beta_{T-1} a_{T-1} \Delta \tilde{P}_T\right) - c_{T-1}(a_{T-1}) &= \sum_{0 \le t < T-1} \left[\Gamma_{t+1} - \beta_t A_t \Delta \tilde{P}_{t+1} - U_t^a(\Gamma_{t+1}) + c_t(A_t)\right] \\
& \quad R + [c_{T-1}(A_{T-1}) - c_{T-1}(a_{T-1})] - U_{T-1}^a(\Gamma_T) \\
& \quad + U_{T-1}^a\left(\Gamma_T - \beta_{T-1} A_{T-1} \Delta \tilde{P}_T + \beta_{T-1} a_{T-1} \Delta \tilde{P}_T\right).
\end{aligned}
$$

By definition of $\mathbb{C}(\beta)$ the sum of the last terms is smaller or equal than 0, and exactly zero when $a_{T-1} = A_{T-1}$. Therefore

$$
\operatorname*{ess\,sup}_{a_{T-1}} \left\{ U_{T-1}^a\left(\Theta + \beta_{T-1} a_{T-1} \Delta \tilde{P}_T\right) - c_{T-1}(a_{T-1}) \right\} = R + \sum_{0 \le t < T-1} \left[\Gamma_{t+1} - \beta_t A_t \Delta \tilde{P}_{t+1} + c_t(A_t) - U_t^a(\Gamma_{t+1})\right].
$$

This shows that at time $T-1$ the Agent chooses $A_{T-1}$ when presented with $(\Theta, \beta)$. If we define $H_T = \Theta$, we are thus entitled to call $H_T$ the value (left hand side or right one) in the above equality. By using backwards induction, as we have often done and hence omit, we have proven that the contract $S$ (defined from $(\Theta, \beta)$) is optimal for the Principal and incentive compatible (notice that automatically $H_0 = R$), and the Agent indeed chooses $A$ under this contract. $\qquad \square$

## REFERENCES

[1] P. Barrieu and N. El Karoui. Inf-convolution of risk measures and optimal risk transfer. *Finance Stoch.*, 9(2):269–298, 2005.

[2] Aharon Ben-Tal and Marc Teboulle. An old-new concept of convex risk measures: the optimized certainty equivalent. *Math. Finance*, 17(3):449–476, 2007.

[3] B. Biais, T. Mariotti, J.C. Rochet, and S. Villeneuve. Large risks, limited liability, and dynamic moral hazard. *Econometrica*, 78, 2010.

[4] J. F. Bonnans and A. Shapiro. Optimization problems with perturbations: a guided tour. *SIAM Rev.*, 40(2):228–264, 1998.

[5] K. Borch. Reciprocal reinsurance treaties seen as a two-person co-operative game. *Skand. Aktuarietidskr.*, 1960:29–58 (1961), 1960.

[6] A. Cadenillas, J. Cvitanić, and F. Zapatero. Optimal risk-sharing with effort and project choice. *J. Econom. Theory*, 133(1):403–440, 2007.

[7] P. Cheridito, M. Kupper, and N. Vogelpoth. Conditional analysis on $R^d$. Forthcoming in Volume 'Set Optimization and Applications - State of the Art'.

[8] Patrick Cheridito, Ulrich Horst, Michael Kupper, and Traian A. Pirvu. Equilibrium pricing in incomplete markets under translation invariant preferences. *Math. Oper. Res.*, 41(1):174–195, 2016.

[9] A. Cherny and M. Kupper. Divergence utilities. Preprint, 2007.

[10] J. Cvitanić, D. Possamaï, and N. Touzi. Moral hazard in dynamic risk management. 2014.

[11] J. Cvitanić and J. Zhang. *Contract theory in continuous-time models*. Springer Finance. Springer, Heidelberg, 2013.

[12] F. Delbaen and W. Schachermayer. A general version of the fundamental theorem of asset pricing. *Math. Ann.*, 300(3):463–520, 1994.

[13] L. Epstein and T. Wang. Intertemporal asset pricing under knightian uncertainty. *Econometrica*, 62(2):pp. 283–322, 1994.

[14] D. Filipović, M. Kupper, and N. Vogelpoth. Separation and duality in locally $L^0$-convex modules. *J. Funct. Anal.*, 256(12):3996–4029, 2009.

[15] D. Filipović, M. Kupper, and N. Vogelpoth. Approaches to conditional risk. *SIAM J. Financial Math.*, 3(1):402–432, 2012.

[16] I. Gilboa and D. Schmeidler. Maxmin expected utility with nonunique prior. *J. Math. Econom.*, 18(2):141–153, 1989.

[17] B. Holmström. Moral hazard and observability. *The Bel Journal of Economics*, 10:74–91, 1979.

[18] B. Holmström and P. Milgrom. Aggregation and linearity in the provision of intertemporal incentives. *Econometrica*, 55(2):303–328, 1987.

[19] F. Maccheroni, M. Marinacci, and A. Rustichini. Ambiguity aversion, robustness, and the variational representation of preferences. *Econometrica*, 74(6):1447–1498, 2006.

[20] H. Ou-Yang. Optimal contracts in a continuous-time delegated portfolio management problem. *Rev. Financ. Stud.*, 16-1:173–208, 2003.

[21] R.T. Rockafellar and S. Uryasev. Conditional value-at-risk for general loss distributions. *Journal of Banking and Finance*, pages 1443–1471, 2002.

[22] Y. Sannikov. A continuous-time version of the principal-agent problem. *Rev. Econom. Stud.*, 75(3):957–984, 2008.

[23] H. Schättler and J. Sung. The first-order approach to the continuous-time principal-agent problem with exponential utility. *J. Econom. Theory*, 61(2):331–371, 1993.

[24] S. Spear and S. Srivastava. On repeated moral hazard with discounting. *Rev. Econom. Stud.*, 54(4):599–617, 1987.

[25] Noah Williams. A solvable continuous time dynamic principal-agent model. *J. Econom. Theory*, 159(part B):989–1015, 2015.

# Measuring the Agency Cost of Debt

## ANTONIO S. MELLO and JOHN E. PARSONS*

**ABSTRACT**

We adapt a contingent claims model of the firm to reflect the incentive effects of the capital structure and thereby to measure the agency costs of debt. An underlying model of the firm and the stochastic features of its product market are analyzed and an optimal operating policy is chosen. We identify the change in operating policy created by leverage and value this change. The model determines the value of the firm and its associated liabilities incorporating the agency consequences of debt.

THE OPTIMAL CAPITAL STRUCTURE for a firm is now widely regarded to be determined by a broad range of factors including a mix of tax effects, the various agency problems associated with different securities, and the various costs of issuing securities, including the costs created by adverse selection. While the existence of a theoretical optimum has been demonstrated in a variety of papers, a less explored area has been the construction of detailed models that enable us to measure each of the relevant factors for a particular company and thereby to determine the actual optimal mix for that firm. This gap in our understanding is particularly glaring in the case of agency costs. In order to allow a careful modelling of strategic relations, the parameters of most agency models are either so simplified that it is impossible to associate them with measurable parameters of a real world case, or else the models simply abstract from certain critical factors—such as a robust measure of price risk—that must be incorporated into any real application. For example, although we now understand that sinking funds, dividend restrictions, and other bond covenants help to resolve the conflict of interest between bond-holders and equity, we do not yet have much in the way of models with which to determine the optimal parameters of these very covenants.

Contingent claims models can provide a consistent framework for multi-period valuation that properly accounts for risk, but they usually abstract from the agency factors entering capital structure decisions. When using a contingent claims model to value a firm's securities it is common to take the value of the firm itself as governed by an exogenously defined stochastic process. The value of the firm's securities are then derived from this underlying value, and, as Merton (1974) points out in his paper on the pricing of risky debt, the Modigliani-Miller theorem obtains so that the value of the firm is independent of the value and the type of debt.

In order to apply the contingent claims techniques to a setting in which agency problems are central, some adaptation is necessary. In this paper we

* Mello is from The Banco de Portugal and Parsons is from Baruch College, C.U.N.Y.

1887

make the value of the firm an endogenous function of (1) an underlying stochastic variable describing the firm's product market and of (2) the management's choice of operating and investment decisions. The management maximizes the value of the equity claim as valued using a traditional contingent claims model, and this in turn determines the actual realized stochastic process that will describe the value of the firm and its debt. Different assumed financial structures will induce different operating strategies and therefore different realized stochastic processes for the value of the firm. The divergence of the chosen operating policy away from the first best operating policy gives rise to an agency cost of debt, and we are able to use the contingent claims model to precisely measure this cost and to identify how it varies with the underlying parameters of the model and with the relative profitability of the firm.

Earlier work on this type of problem includes Brennan and Schwartz (1984) and Fischer, Heinkel, and Zechner (1989). The former authors analyze the equity owners' optimal reinvestment and external financing policy over time given constraints imposed by pre-existing bond covenants. The latter consider the equity owners' optimal recapitalization policy in light of transaction costs and the tax benefits of debt, and they are able to then prescribe optimal ex ante call values to be included in the debt contract as it is originally written.

In this paper we focus on a specific production problem and analyze how the existence of debt directly changes the equity owners' choice of an operating policy for the business. The Brennan and Schwartz (1985) contingent claims model for valuing a mine is extended to incorporate the financial structure and to recognize the effects of the agency problems. An interesting benefit of using their model is that we are able to identify precisely the changes in the operating policy of the mine that are induced by the outstanding debt and thereby to directly relate the agency costs of debt with clearly suboptimal decisions in real production. The agency problems that arise are the underinvestment problem identified by Myers (1977), as well as the costs of bankruptcy. Our model measures these costs and thereby compares the tax benefits of debt with the agency costs of debt.

In Section I we extend the standard contingent claims model of the firm to incorporate the incentive effects of leverage on the firm's choice of operating policy. In Section II we apply this model to the task of measuring the agency costs. In Section III we present some numerical results and perform some analysis. We conclude in Section IV.

## I. A Contingent Claims Valuation of a Mine in the Presence of Agency Costs

Brennan and Schwartz (1985) analyze a firm that owns a mine with a commodity inventory, $Q$. When the mine is open the commodity is extracted at a constant annual rate, $q$, and at a constant real average annual cost, $a$. When the mine is closed a constant real annual maintenance cost, $m$, is

incurred. Corporate taxes are paid at rate $\tau$ on net income, and it is assumed that full offsets are allowed. At any point in time the mine can be closed at a real cost $k_1$ and reopened at a real cost $k_2$. The mine can also be costlessly abandoned.

Several crucial assumptions are made on the stochastic structure of the commodity price. First, the real spot price of the commodity, $s$, is determined in a competitive market and follows the exogenous process

$$ds = \mu s dt + \sigma s dz, \qquad (1)$$

where $dz$ is the increment to a standard Gauss-Wiener process; $\sigma$, the instantaneous standard deviation of the spot price, is assumed to be known and constant; and $\mu$ is the instantaneous drift in the real price. Second, it is assumed that there is a traded futures contract on the commodity. Then, following Ross (1978), if the convenience yield on the commodity is a constant proportion of the spot price, $\kappa(s) = \kappa s$, and if there exists a known constant real interest rate, $r$, the real price of a futures contract maturing in $\bar{t}$ periods is given by $f(s, \bar{t}) = se^{(r-\kappa)\bar{t}}$.

The market value of the mine, $v \equiv v(s, Q; j, \phi)$, is a function of the current commodity price, $s$, of the inventory, $Q$, of whether the mine is currently closed or open, $j = 1, 2$, and of the optimal operating policy, $\phi$. An operating policy is described by three functions defining three critical commodity prices: $s_0(Q)$, the price, for a given inventory level, at which the mine is abandoned if it is already closed, $s_1(Q)$, the price for a given inventory level, at which the mine is closed if it was previously open, and $s_2(Q)$, the price, for a given inventory level, at which the mine is opened if it was previously closed, $\phi = (s_0, s_1, s_2)$. Throughout the remainder of the paper we suppress the argument $Q$, and write each of these functions simply as $s_i$, $i = 0, 1, 2$. The extraction rate for an open mine is assumed constant at $q$. Applying Ito's lemma of stochastic calculus the instantaneous change in the value of the mine is given by $dv = v_s ds + v_Q dQ + {}^1\!/_2 v_{ss}(ds)^2$. The cash flow from the mine is $[q(s-a)(j-1) - m(2-j)](1-\tau)$. Using an arbitrage argument similar to Black and Scholes (1973) the differential equation governing the value of the closed mine is written

$$^1\!/_2 \sigma^2 s^2 v_{ss}(s, Q; 1) + (r - \kappa) s v_s(s, Q; 1) - m(1 - \tau) - rv(s, Q; 1) = 0, \quad (2)$$

and the value of the open mine is written

$$^1\!/_2 \sigma^2 s^2 v_{ss}(s, Q; 2) + (r - \kappa) s v_s(s, Q; 2) - q v_Q(s, Q; 2)$$
$$+ q(s - a)(1 - \tau) - rv(s, Q; 2) = 0. \quad (3)$$

Associated with this pair of equations are four boundary conditions:

$$v(s, 0; j) = 0, \qquad (4)$$

$$v(s_0, Q; 1) = 0, \qquad (5)$$

$$v(s_1, Q; 2) = \max\{v(s_1, Q; 1) - k_1(1 - \tau), 0\}, \qquad (6)$$

$$v(s_2, Q; 1) = v(s_2, Q; 2) - k_2(1 - \tau). \qquad (7)$$

The first best operating policy $\phi^{FB} \equiv (s_0^{FB}, s_1^{FB}, s_2^{FB})$ is characterized by the following first order conditions:

$$v_s(s_0^{FB}, Q; 1) = 0, \tag{8}$$

$$v_s(s_1^{FB}, Q; 2) = \begin{cases} v_s(s_1^{FB}, Q; 1) & \text{if } v(s_1^{FB}, Q; 1) - k_1(1 - \tau) \geq 0 \\ 0 & \text{if } v(s_1^{FB}, Q; 1) - k_1(1 - \tau) < 0, \end{cases} \tag{9}$$

$$v_s(s_2^{FB}, Q; 1) = v_s(s_2^{FB}, Q; 2). \tag{10}$$

solving equations (2) and (3) subject to boundary conditions (4)–(10), we derive simultaneously the first best value of the mine and the first best operating policy, $v^{FB}$ and $\phi^{FB}$.

If the firm is financed in part with debt, then the first best solution will not generally be chosen by managers acting in the interest of equity holders because the debt creates agency problems in the operation of the mine. The actual value of the firm will be determined by the operating policy chosen to maximize the value of the levered equity, the second best operating policy. To correctly value the firm and its associated liabilities we incorporate the effect of leverage into the Brennan and Schwartz (1985) model's derivation of the optimal operating policy.

We assume that the mine is financed in part with a bond described by a time path of the outstanding principal balance $P(t)$ and by a constant continuous coupon rate $c$. The interest payments on the bond, $cP(t)$, are tax deductible. We assume that there will be some point in time, $T$, such that $\forall t \geq T$, $P(t) = 0$, and we will call this the final maturity date of the bond. This covers a large range of possible debt structures. For example, a bond with constant amortization would satisfy the condition that $\forall t \leq T$, $P_t + cP = \delta$. A bond with a balloon payment at maturity can be approximated by a bond with zero principal payments until close to maturity, $\forall t < T - \varepsilon$, $P_t = 0$ and with continuous and quickly increasing principal payments as maturity approaches, $t > T - \varepsilon$ then $P_t \to \infty$ as $t \to T$. If $T$ is much larger than the life of the mine, then this bond is comparable to a perpetuity.[1] When the assumed bond structure is very complicated it may have the appearance of a debt policy rather than of a single instrument. Our model values the firm given any assumed structure or policy. However, we do not allow the structure to be costlessly altered ex post to avoid bankruptcy. We solve for the optimal debt policy ex ante.

Before going back to revalue the levered mine it is necessary to value the equity since it is this value that will be maximized in choosing the firm's actual operating policy. The market value of the equity, $e \equiv e(s, Q, t; j, \phi')$, is a function of the current commodity price, $s$, of the inventory, $Q$, of the current time period, $t$, of whether the mine is currently closed or open, $j = 1, 2$, and of the modified operating policy, $\phi'$. The modified operating

---

[1] For other payment structures that can be solved using the methodology of this paper see Mello and Parsons (1991).

policy acknowledges the right of the equity owners to default on the bond and is described by three functions defining three critical commodity prices: $s_d(Q, t)$ is the price, for any given inventory at time $t$, at which the equity owners default, while $s_1(Q, t)$ and $s_2(Q, t)$ are, as before, the prices, for any given inventory at time $t$, at which the mine is closed or opened, respectively, $\phi' = (s_d, s_1, s_2)$. Once again, we generally suppress the arguments $Q$ and $t$, and write each of these functions simply as $s_i$, $i = d, 1, 2$. Upon default the firm is sold at its then current value, $v(s, Q, t; j)$ and these proceeds go to the bondholders. This possibility of default, of course, gives the equity owners a compound call option on the value of the mine. Consequently the value of the equity is time dependent and hence path dependent. Again, applying Ito's lemma, the instantaneous change in the value of the equity is given by $de = e_s ds + e_Q dQ + e_t dt + {}^1\!/_2 e_{ss}(ds)^2$. The cash flow from the equity is $[q(s - a)(j - 1) - m(2 - j)](1 - \tau) + P_t - (1 - \tau)cP(t)$. The last two terms are the principal payment and the after tax interest payment on the bond. Then the differential equation governing the value of the equity when the mine is closed is:

$$
{}^1\!/_2 \sigma^2 s^2 e_{ss}(s, Q, t; 1) + (r - \kappa) s e_s(s, Q, t; 1) + e_t(s, Q, t; 1)
$$
$$
- m(1 - \tau) + P_t - (1 - \tau)cP(t) - re(s, Q, t; 1) = 0, \quad (11)
$$

and when the mine is open is:

$$
{}^1\!/_2 \sigma^2 s^2 e_{ss}(s, Q, t; s) + (r - \kappa) s e_s(s, Q, t; 2) - q e_Q(s, Q, t; 2)
$$
$$
+ e_t(s, Q, t; 2) + q(s - a)(1 - \tau) + P_t - (1 - \tau)cP(t)
$$
$$
- re(s, Q, t; 2) = 0. \quad (12)
$$

Associated with this pair of equations are also four boundary conditions:

$$
e(s, 0, t; j) = 0, \quad (13)
$$
$$
e(s_d, Q, t; 1) = 0, \quad (14)
$$
$$
e(s_1, Q, t; 2) = \max\{e(s_1, Q, t; 1) - k_1(1 - \tau), 0\}, \quad (15)
$$
$$
e(s_2, Q, t; 1) = e(s_2, Q, t; 2) - k_2(1 - \tau). \quad (16)
$$

The equity owner's optimal operating policy, $\phi'^P = (s_d^P, s_1^P, s_2^P)$ is characterized by the following first order conditions:

$$
e_s(s_d^P, Q, t; 1) = 0, \quad (17)
$$

$$
e_s(s_1^P, Q, t; 2) = \begin{cases} e_s(s_1^P, Q, t; 1) & \text{if } e(s_1^P, Q, t; 1) - k_1(1 - \tau) \geq 0 \\ 0 & \text{if } e(s_1^P, Q, t; 1) - k_1(1 - \tau) < 0, \end{cases} \quad (18)
$$

$$
e_s(s_2^P, Q, t; 1) = e_s(s_2^P, Q, t; 2), \quad (19)
$$

The value of equity, $e^P$, and the optimal operating policy, $\phi'^P = (s_d^P, s_1^P, s_2^P)$, are derived simultaneously as the solution to the two differential equations (11) and (12) using boundary conditions (13)–(19).

To determine the value of the levered firm it is necessary to solve the pair of differential equations:

$$\tfrac{1}{2}\sigma^2 s^2 v_{ss}(s,Q,t;1) + (r-\kappa)s v_s(s,Q,t;1) + v_t(s,Q,t;1)$$
$$- m(1-\tau) + \tau c P(t) - r v(s,Q,t;1) = 0, \quad (20)$$

and

$$\tfrac{1}{2}\sigma^2 s^2 v_{ss}(s,Q,t;2) + (r-\kappa)s v_s(s,Q,t;2) - q v_Q(s,Q,t;2)$$
$$+ v_t(s,Q,t;2) + q(s-a)(1-\tau) + \tau c P(t) - r v(s,Q,t;2) = 0. \quad (21)$$

along with boundary conditions based upon the operating policy that is optimal for the equity owners:

$$v(s,0,t;j) = 0, \quad (22)$$

$$v(s,Q,t;j) = v^{FB}(s,Q;j), \qquad \forall t \geq T \quad (23)$$

$$v(s_d^P, Q, t; 1) = \alpha v^{FB}(s_d^P, Q; 1), \quad (24)$$

$$v(s_1^P, Q, t; 2) = \max\{v(s_1^P, Q, t; 1) - k_1(1-\tau), 0\}, \quad (25)$$

$$v(s_2^P, Q, t; 1) = v(s_2^P, Q, t; 2) - k_2(1-\tau). \quad (26)$$

The value for the levered mine calculated using this system of equation is denoted $v^P$.

Boundary condition (24) requires some comment. Upon default the firm is put to the bondholder. The case in which the firm is subsequently operated according to the first best operating policy is equivalent to setting $\alpha = 1$. Another more general case incorporates the possibilities that either (1) there are costs of financial distress associated with bankruptcy, or (2) the bondholder cannot operate the firm and must reorganize it with a similar debt/equity structure—thereby reproducing the agency problem. This case is described by letting $\alpha \in [0,1)$. The parameter $\alpha$ measures the significance of the costs of financial distress, and as $\alpha$ approaches zero these agency costs increase.

The value for the outstanding bond is the difference between the total value of the mine and the value of the equity:

$$b^P = v^P - e^P. \quad (27)$$

To illustrate the model we calculated values for $v^P$, $e^P$, and $b^P$ for a hypothetical example with input parameter listed in Table I. To our knowledge there is no closed-form solution to these various systems of equations. It is, however, possible to solve this system of equations using numerical methods as we have done for the hypothetical mine. The input parameters for our example are given in Table I. The critical commodity prices characterizing the equity owners' optimal operating policy at the initial inventory and at $t = 0$, $s_d^P(Q,t)$, $s_1^P(Q,t)$, and $s_2^P(Q,t)$, are displayed in Table II and contrasted with the critical commodity prices characterizing the first best operat-

<div align="center">

Table I

## Data for the Hypothetical Mining Firm

</div>

| | |
|---|---|
| Total inventory in the ground: | $Q = 150$ million pounds |
| Annual real production for an open mine: | $q = 10$ million pounds |
| Average real production costs: | $a = \$0.50/\text{pound}$ |
| Maintenance costs for a closed mine: | $m = \$0.0/\text{year}$ |
| Closing and opening costs: | $k_1 = k_2 = \$2$ million |
| Real interest rate: | $r = 2\%$ annually |
| Commodity price variance: | $\sigma^2 = 8\%$ annually |
| Convenience yield: | $\kappa = 1.5\%$ annually |
| Corporate income tax rate: | $\tau = 34\%$ |

<div align="center">

Table II

## The Levered Firm's Choice of an Operating Policy

</div>

| | |
|---|---|
| The Bond Contract | |
| par value | \$5.24 million |
| coupon rate | 2% |
| annual debt service | \$0.4 million |
| final maturity date | 15 years |
| Factor of firm's first best value at bankruptcy: | $\alpha = 0$ |

| Critical Commodity Prices[a] ($/pound) | First Best Operating Policy | Equity Owner's Optimal Operating Policy |
|---|---|---|
| abandonment/default | $s_0^{FB} = 0.00$ | $s_d^P = 0.40$ |
| closing | $s_1^{FB} = 0.59$ | $s_1^P = 0.54$ |
| opening | $s_2^{FB} = 0.84$ | $s_2^P = 0.79$ |

[a] All values for the critical commodity prices in the optimal operating policies are calculated at the initial inventory given in Table I, 150 million pounds, and at $t = 0$.

ing policy at the initial inventory, $s_0^{FB}(Q)$, $s_1^{FB}(Q)$, and $s_2^{FB}(Q)$. The values for the levered firm, $v^P$, levered equity, $e^P$, and for the bond, $b^P$, are displayed in Table III.

## II. Measuring the Agency Cost of Debt

It is important to note that in general the operating policy chosen to maximize the value of the equity claim will not be identical with the first best operating policy, $(s_d^P, s_1^P, s_2^P) \neq (s_0^{FB}, s_1^{FB}, s_2^{FB})$.[2] Without any agency costs of debt the value of the levered firm would be the first best value of the firm plus the interest tax shield of debt. Each added unit of debt increases the value of the firm by the value of its associated interest tax shields. The presence of agency costs modifies this. At first, with no debt outstanding, a

[2] More accurately, $\forall Q > 0$, $t < T$, $(s_d^P(Q, t), s_1^P(Q, t), s_2^P(Q, t)) \neq (s_0^{FB}(Q), s_1^{FB}(Q), s_2^{FB}(Q))$.

<div style="text-align:center">

**Table III**

**The Value of the Levered Firm and Its Liabilities**

</div>

| The Bond Contract | |
| --- | --- |
| par value | $5.24 million |
| coupon rate | 2% |
| annual debt service | $0.4 million |
| final maturity date | 15 years |

Factor of firm's first best value at bankruptcy: $\alpha = 0$

| Commodity Price $s$ | Firm Value $v^P(s, Q, t; j)$ | | Equity Value $e^P(s, Q, t; j)$ | | Bond Value $b^P(s, Q, t; j)$ | |
| --- | --- | --- | --- | --- | --- | --- |
| | (closed) $j = 1$ | (open) $j = 2$ | (closed) $j = 1$ | (open) $j = 2$ | (closed) $j = 1$ | (open) $j = 2$ |
| 0.05 | 0 | | 0 | | 0 | |
| 0.10 | 0 | | 0 | | 0 | |
| 0.15 | 0 | | 0 | | 0 | |
| 0.20 | 0 | | 0 | | 0 | |
| 0.25 | 0 | | 0 | | 0 | |
| 0.30 | 0 | | 0 | | 0 | |
| 0.35 | 0 | | 0 | | 0 | |
| 0.40 | 0 | | 0 | | 0 | |
| 0.45 | 2.71 | | 0.32 | | 2.39 | |
| 0.50 | 5.80 | | 1.84 | | 3.97 | |
| 0.55 | 8.80 | 6.90 | 4.14 | 3.45 | 4.66 | 3.46 |
| 0.60 | 12.08 | 10.79 | 7.01 | 6.83 | 5.07 | 3.96 |
| 0.65 | 15.44 | 15.26 | 10.30 | 10.69 | 5.14 | 4.57 |
| 0.70 | 19.15 | 19.58 | 13.93 | 14.81 | 5.22 | 4.73 |
| 0.75 | 23.03 | 24.07 | 17.97 | 19.08 | 5.16 | 4.99 |
| 0.80 | | 28.52 | | 23.44 | | 5.09 |
| 0.85 | | 32.99 | | 27.84 | | 5.15 |
| 0.90 | | 37.46 | | 32.28 | | 5.18 |
| 0.95 | | 41.92 | | 36.72 | | 5.20 |
| 1.00 | | 46.40 | | 41.18 | | 5.22 |

marginal increase in debt has the same effect as before, increasing the value of the firm above the first best value by the size of the interest tax shields. As the size of debt increases, however, the marginal agency costs grow so that for large values of debt the total agency costs may far outweigh the total tax shields making the value of the levered firm less than the first best.

Our objective is to directly measure the agency costs associated with a particular financial structure. To do so we need to separate the effect of the tax shield for any outstanding bond: we define $\eta$ as the value of the interest tax shield earned by the firm operated according to the policy $\phi'$. The value of the interest tax shield is the solution to the following pair of differential equations:

$$\tfrac{1}{2}\sigma^2 s^2 \eta_{ss}(s, Q, t; 1) + (r - \kappa)s\eta_s(s, Q, t; 1) + \eta_t(s, Q, t; 1)$$
$$+ \tau c P(t) - r\eta(s, Q, t; 1) = 0, \quad (28)$$

and,

$$\tfrac{1}{2}\sigma^2 s^2 \eta_{ss}(s,Q,t;2) + (r-\kappa)s\eta_s(s,Q,t;2) - q\eta_Q(s,Q,t;2)$$
$$+ \eta_t(s,Q,t;2) + \tau cP(t) - r\eta(s,Q,t;2) = 0. \quad (29)$$

along with the boundary conditions:

$$\eta(s,0,t;j) = 0, \quad (30)$$
$$\eta(s,Q,t;j) = 0, \quad \forall t \geq T \quad (31)$$
$$\eta(s_d^P,Q,t;1) = 0, \quad (32)$$
$$\eta(s_1^P,Q,t;2) = \eta(s_1^P,Q,t;1), \quad (33)$$
$$\eta(s_2^P,Q,t;1) = \eta(s_2^P,Q,t;2). \quad (34)$$

The agency costs of the debt $P(t)$ can then be defined as

$$\psi^P(s,Q,t;j) \equiv v^{FB}(s,Q;j) - [v^P(s,Q,t;j) - \eta^P(s,Q,t;j)]. \quad (35)$$

This is a precise measure of the value lost when the equity owners, because of the outstanding debt, change the operating policy from the first best.

Since the operating policy chosen to maximize the value of the equity is not the first best operating policy, the value of the levered firm is less than the first best value of the firm plus the interest tax shields, $v^P < v^{FB} + \eta^P$, the difference being the agency cost of debt. In Table IV the values for $v^P$ are compared against the values for $v^{FB}$ for the sample parameters described above, and the interest tax shields, $\eta^P$, and the agency costs of debt, $\psi^P$, are calculated.

The total agency costs reported in the tables are the consequence of three different changes in the firm's operating policy. First, the abandonment decision and its consequences is a straightforward example of the deadweight costs of financial distress that have been much discussed in the literature (see Shapiro and Titman (1986)). Second, the shareholders change the firm's policy for when to close the mine, $s_1^P < s_1^{FB}$, keeping it open longer than is Pareto optimal in the face of a falling commodity price. By spending the money to close the mine the firm saves on operating costs and preserves the limited inventory until the price rises again. However, the shareholders will bear the full expense of closure and do not enjoy the full benefits. They gamble that the commodity price may rise again without having fallen too far. In this case they will have avoided paying out of their own pockets the fixed cost of closure and then once again the fixed cost of reopening the mine. Third, the shareholders are also changing the firm's policy for when to reopen the mine, $s_2^P < s_2^{FB}$, opening it sooner than is Pareto optimal in the face of a rising commodity price. The shareholders have an interest in extracting the inventory as quickly as possible, since in the event of a future price decline they may have to put to the debtholders whatever remains of the inventory.

## Table IV
## The Agency Cost of Debt

The Bond Contract

| | |
|---|---|
| par value | $5.24 million |
| coupon rate | 2% |
| annual debt service | $0.4 million |
| final maturity date | 15 years |
| Factor of firm's first best value at bankruptcy: | $\alpha = 0$ |

| Commodity Price $s$ | First Best Firm Value $v^{FB}(s, Q, t; j)$ | | Levered Firm Value $v^P(s, Q, t; j)$ | | Tax Shields $\eta^P(s, Q, t; j)$ | | Agency Costs | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | Absolute Value $\psi^P \equiv v^{FB} + \eta^P - v^P$ | | % of First Best $\psi^P/v^{FB}$ | |
| | (closed) $j=1$ | (open) $j=2$ | (closed) $j=1$ | (open) $j=2$ | (closed) $j=1$ | (open) $j=2$ | (closed) $j=1$ | (open) $j=2$ | (closed) $j=1$ | (open) $j=2$ |
| 0.05 | 0.00 | | 0 | | 0 | | 0.00 | | 100 | |
| 0.10 | 0.05 | | 0 | | 0 | | 0.05 | | 100 | |
| 0.15 | 0.25 | | 0 | | 0 | | 0.25 | | 100 | |
| 0.20 | 0.69 | | 0 | | 0 | | 0.69 | | 100 | |
| 0.25 | 1.39 | | 0 | | 0 | | 1.39 | | 100 | |
| 0.30 | 2.36 | | 0 | | 0 | | 2.36 | | 100 | |
| 0.35 | 3.60 | | 0 | | 0 | | 3.60 | | 100 | |
| 0.40 | 5.11 | | 0 | | 0 | | 5.11 | | 100 | |
| 0.45 | 6.90 | | 2.71 | | 0.136 | | 4.33 | | 62.7 | |
| 0.50 | 8.97 | | 5.80 | | 0.212 | | 3.38 | | 37.6 | |
| 0.55 | 11.31 | | 8.80 | 6.90 | 0.244 | 0.244 | 2.75 | 3.33 | 24.3 | 33.3 |
| 0.60 | 13.94 | 12.67 | 12.08 | 10.79 | 0.257 | 0.257 | 2.12 | 2.13 | 15.2 | 16.8 |
| 0.65 | 16.86 | 16.14 | 15.44 | 15.26 | 0.262 | 0.261 | 1.68 | 1.14 | 10.0 | 7.1 |
| 0.70 | 20.07 | 20.04 | 19.15 | 19.58 | 0.264 | 0.264 | 1.18 | 0.73 | 5.9 | 3.6 |
| 0.75 | 23.59 | 24.18 | 23.03 | 24.07 | 0.266 | 0.265 | 0.82 | 0.38 | 3.5 | 1.6 |
| 0.80 | 27.40 | 28.47 | | 28.52 | | 0.266 | 0.46 | 0.22 | 1.7 | 0.8 |
| 0.85 | | 32.84 | | 32.99 | | 0.266 | | 0.12 | | 0.4 |
| 0.90 | | 37.26 | | 37.46 | | 0.266 | | 0.07 | | 0.2 |
| 0.95 | | 41.70 | | 41.92 | | 0.266 | | 0.05 | | 0.1 |
| 1.00 | | 46.15 | | 46.40 | | 0.266 | | 0.02 | | 0.0 |

## III. Results

How significant are the agency costs of debt? Suppose that the current price of the commodity is $0.80/pound. An open mine would have annual revenues at this price of $8 million, annual costs of $5 million, and a net cash flow after tax of $1.98 million. If the firm operating an open mine has outstanding a bond with 15-year maturity and constant annual debt service payments of $0.4 million we can see from Table III that its present value is $28.52 million. The bond would have a market value of $5.09 million, that is less than 18% of the firm value, a very low debt-to-value ratio. Moreover, the firm's current annual cash flow is five times its annual debt obligation. Clearly the probability of bankruptcy appears very small, and many would imagine that the agency costs of the debt should be correspondingly miniscule. From Table IV we read that the agency costs of this quantity of debt are $0.22 million, or eight-tenths of a percent of firm value. In terms of the amount of debt sold, however, these agency costs are close to 4.3%, a very large value. This should be compared to other costs such as underwriting fees and administrative expenses which are usually 1.3% of the value of a debt offering according to Mikkelson and Partch (1986).

As mentioned earlier, this total agency cost is a combination of various factors—the suboptimal opening and closure policies and the dead weight cost of bankruptcy. To disentangle these causes and to explore the significance of the pure operational factors we reparameterize our model setting $\alpha = 1$ so that there are absolutely no dead weight costs associated with bankruptcy: the bondholders receive the first best value of the firm. The results for this case are displayed in Table V. When, as before, the current price of the commodity is $0.80/pound the firm operating an open mine and with outstanding a bond with 15-year maturity and constant annual debt service payments of $0.4 million has a present value of $28.64 million. Not shown in the table, the bond would have a market value of $5.20 million, again close to 18% of the firm value. The agency costs of debt in this case are $0.10 million, about one half of the total agency costs from the previous example. The agency costs of debt amount to three-tenths of a percent of firm value or almost 2% of the value of debt sold.

Agency costs of this magnitude would certainly be an important determinant of the firm's capital structure decision even though the firm appears far from bankrupt. Table VI contains results also displayed in Figures 1 and 2. In Figure 1 we graph the levered firm's value as a function of its debt-to-value ratio. When the commodity price is $0.65/pound and the mine is open the first best value of the firm is $16.141 million. An outstanding 15-year bond with a 2% coupon and constant annual debt service payments totalling $0.01 million offers tax shields worth $0.007 million. The agency costs are $0.005 million and the value of the levered firm is $16.143 million. The bond is valued at $0.228 million giving a debt-to-value ratio of less than 2%. Higher debt loads only lower the value of the firm since the marginal tax shields are less than the marginal agency costs of debt. When the commodity price is

## Table V
## The Agency Cost of Debt When Bankruptcy is Costless

The Bond Contract

| | |
|---|---|
| par value | $5.24 million |
| coupon rate | 2% |
| annual debt service | $0.4 million |
| final maturity date | 15 years |

Factor of firm's first best value at bankruptcy:  $\alpha = 1$

| Commodity Price $s$ | First Best Firm Value $v^{FB}(s,Q,t;l)$ | | Levered Firm Value $v^{P}(s,Q,t;j)$ | | Tax Shields $\eta^{P}(s,Q,t;j)$ | | Agency Costs | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | Absolute Value $\psi^{P} \equiv \omega^{FB} + \eta^{P} - v^{P}$ | | % of First Best $\psi^{P}/v^{FB}$ | |
| | (closed) $j=1$ | (open) $j=2$ | (closed) $j=1$ | (open) $j=2$ | (closed) $j=1$ | (open) $j=2$ | (closed) $j=1$ | (open) $j=2$ | (closed) $j=1$ | (open) $j=2$ |
| 0.05 | 0.00 | | —[a] | | | | | | | |
| 0.10 | 0.05 | | — | | | | | | | |
| 0.15 | 0.25 | | — | | | | | | | |
| 0.20 | 0.69 | | — | | | | | | | |
| 0.25 | 1.39 | | — | | | | | | | |
| 0.30 | 2.36 | | — | | | | | | | |
| 0.35 | 3.60 | | — | | | | | | | |
| 0.40 | 5.11 | | — | | | | | | | |
| 0.45 | 6.90 | | 6.58 | | 0.136 | | 0.45 | | 6.6 | |
| 0.50 | 8.97 | | 8.35 | | 0.212 | | 0.82 | | 9.2 | |
| 0.55 | 11.31 | | 10.47 | 8.60 | 0.244 | | 1.08 | 1.63 | 9.5 | 16.3 |
| 0.60 | 13.94 | 12.67 | 13.13 | 12.11 | 0.257 | 0.244 | 1.07 | 0.81 | 7.7 | 6.4 |
| 0.65 | 16.86 | 16.14 | 16.19 | 15.86 | 0.262 | 0.257 | 0.93 | 0.55 | 5.5 | 3.4 |
| 0.70 | 20.07 | 20.04 | 19.60 | 20.03 | 0.264 | 0.261 | 0.73 | 0.27 | 3.6 | 1.4 |
| 0.75 | 23.59 | 24.18 | 23.28 | 24.26 | 0.266 | 0.264 | 0.75 | 0.19 | 2.4 | 0.8 |
| 0.80 | 27.40 | 28.47 | | 28.64 | | 0.265 | | 0.10 | | 0.3 |
| 0.85 | | 32.84 | | 33.06 | | 0.266 | | 0.05 | | 0.2 |
| 0.90 | | 37.26 | | 37.50 | | 0.266 | | 0.03 | | 0.1 |
| 0.95 | | 41.70 | | 41.95 | | 0.266 | | 0.02 | | 0.0 |
| 1.00 | | 46.15 | | 46.41 | | 0.266 | | 0.01 | | 0.0 |

[a] No values are shown for a levered firm when the commodity price is less than $0.45/pound. If the price declines to $0.40/pound, then the levered firm defaults. At this point the firm is reorganized and its value becomes the first best.

## Table VI
## **The Optimal Quantity of Debt**

| The Bond Contract | |
| --- | --- |
| coupon rate | 2% |
| final maturity date | 15 years |
| Factor of firm's first best value at bankruptcy: | $\alpha = 0$ |

### Panel A: Initial Commodity Price, $s = \$0.65/\text{pound}$

| Fixed Annual Bond Payment Principal + Interest $[P(t) - P(t+1)] + cP(t)$ | First Best Firm Value $v^{FB}(s, Q, t; 2)$ | Tax Shields $\eta^P(s, Q, t; 2)$ | Agency Costs $\psi^P(s, Q, t; 2)$ | Levered Firm Value $v^P(s, Q, t; 2)$ |
| --- | --- | --- | --- | --- |
| 0.000 | 16.141 | 0.000 | 0 | 16.141 |
| 0.005 | | 0.003 | 0.001 | 16.143 |
| 0.010 | | 0.007 | 0.005 | 16.143 |
| 0.015 | | 0.010 | 0.013 | 16.138 |
| 0.020 | | 0.013 | 0.024 | 16.130 |
| 0.025 | | 0.017 | 0.030 | 16.128 |
| 0.030 | | 0.020 | 0.035 | 16.126 |
| 0.035 | | 0.023 | 0.041 | 16.123 |
| 0.040 | | 0.027 | 0.074 | 16.094 |
| 0.045 | | 0.030 | 0.081 | 16.090 |
| 0.050 | | 0.033 | 0.107 | 16.067 |
| 0.1 | | 0.067 | 0.395 | 15.813 |
| 0.2 | | 0.133 | 0.787 | 15.487 |
| 0.3 | | 0.197 | 1.142 | 15.196 |
| 0.4 | | 0.261 | 1.144 | 15.258 |
| 0.5 | | 0.322 | 1.470 | 14.993 |
| 0.6 | | 0.378 | 1.466 | 15.053 |
| 0.7 | | 0.432 | 2.012 | 14.561 |
| 0.8 | | 0.486 | 2.046 | 14.581 |
| 0.9 | | 0.487 | 2.494 | 14.134 |
| 1.0 | | 0.503 | 3.038 | 13.606 |

### Panel B: Initial Commodity Price, $s = \$1.00/\text{pound}$

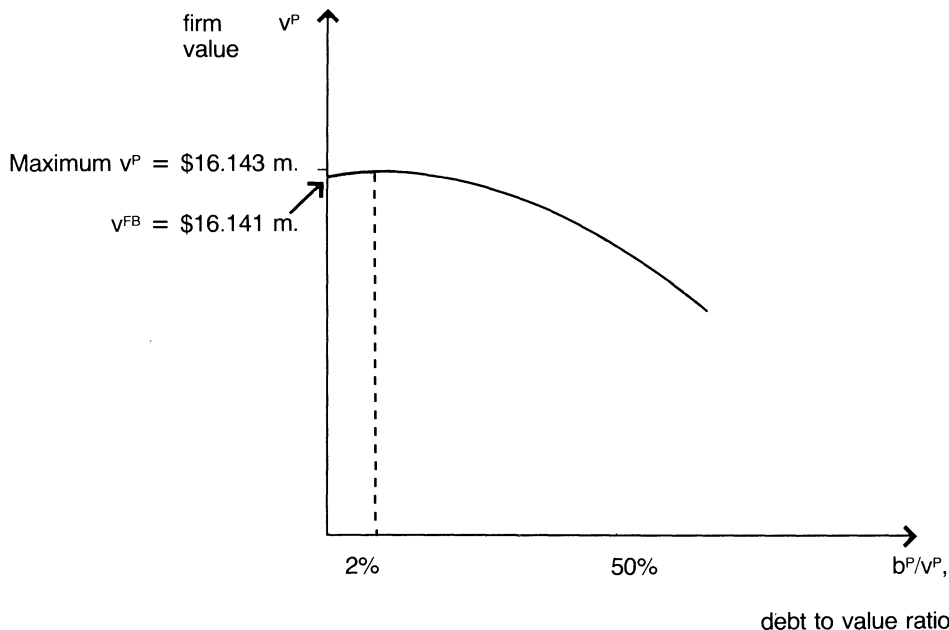| Fixed Annual Bond Payment Principal + Interest $[P(t) - P(t+1)] + cP(t)$ | First Best Firm Value $v^{FB}(s, Q, t; 2)$ | Tax Shields $\eta^P(s, Q, t; j)$ | Agency Costs $\psi^P(s, Q, t; 2)$ | Levered Firm Value $v^P(s, Q, t; 2)$ |
| --- | --- | --- | --- | --- |
| 0.0 | 46.152 | 0.000 | 0 | 46.152 |
| 0.5 | | 0.333 | 0.036 | 46.449 |
| 1.0 | | 0.660 | 0.146 | 46.691 |
| 1.5 | | 0.974 | 0.442 | 46.685 |
| 2.0 | | 1.243 | 1.686 | 45.709 |
| 2.5 | | 1.462 | 2.698 | 44.916 |
| 3.0 | | 1.610 | 6.697 | 29.690 |
| 3.5 | | 1.615 | 10.497 | 30.237 |
| 4.0 | | 1.381 | 20.849 | 23.224 |
| 4.5 | | 1.063 | 27.867 | 18.249 |

**Figure 1. The effect of debt on the value of the firm.** The value of a levered firm owning a mine with parameters specified in Table I. The current commodity price is $0.65/pound. The bond outstanding has a maturity of 15 years and a coupon rate of 2%. As the constant annual debt service payments are increased the debt-to-value ratio increases. The value of the levered firm initially increases due to the marginal benefits of interest tax shields. As the debt-to-value ratio further increases the marginal agency costs rise and the value of the firm begins to fall.

higher, $1.00/pound, the marginal agency costs are lower. This can be seen in Figure 2 where the optimal debt-to-value ratio is much higher. The first best value of the firm is $46.152 million. The interest tax shields associated with a 15-year 2% bond with annual debt service payments of $1.0 million raises the firm value to $46.691 million. The tax shields are valued at $0.66 million and agency costs equal to $0.146 million. The debt to value ratio is now just under 28%. Higher debt ratios lower the value of the firm overall.

It is interesting to explore the consequences of varying the structure of the debt on the total agency costs and the value of the firm. A popular rule of thumb is to match the maturity structure of the firm's liabilities to the maturity structure of its assets. In the case at hand, this rule of thumb would suggest that a constantly amortized bond matching the constant extraction rate of the mine and with a maturity matching the life of the inventory of the mine should be optimal. However, this is not necessarily correct. The maturity structure of the assets is very complicated, and ultimately depends upon the stochastic nature of the commodity price and the operating options available to the firm. These in turn depend upon the capital structure of the firm and the operating policy it induces. To highlight the simultaneity we contrast the results of two numerical examples.
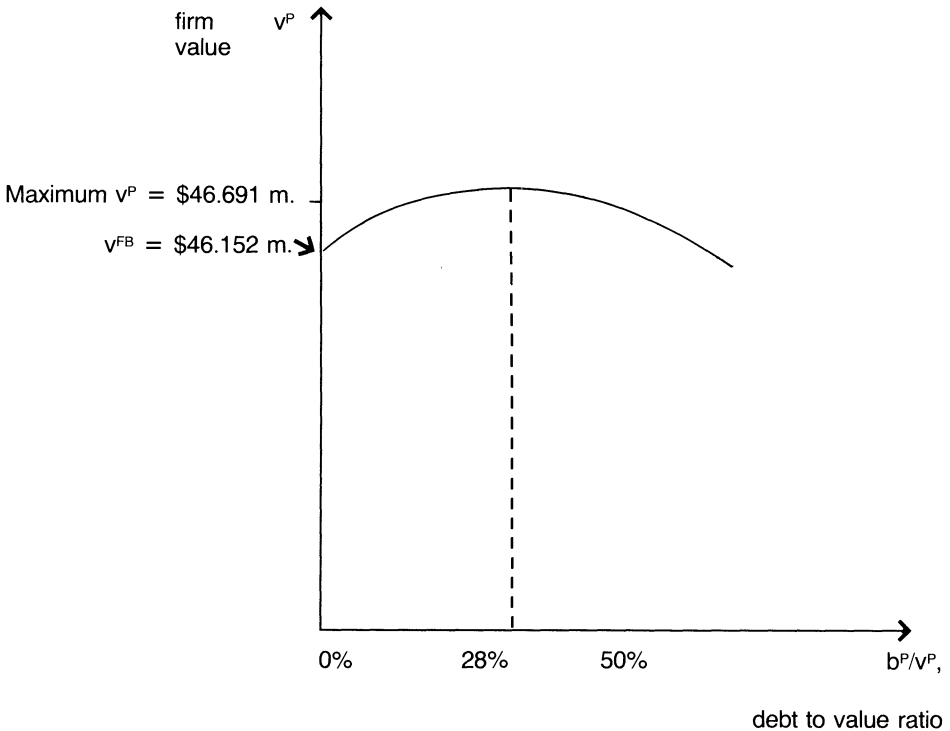
**Figure 2. The effect of debt on the value of the firm at a higher current commodity price, $s$ = \$1.00 / pound.** The value of a levered firm owning a mine with parameters specified in Table I. The bond outstanding has a maturity of 15 years and a coupon rate of 2%. As the constant annual debt service payments are increased the debt-to-value ratio increases. The value of the levered firm initially increases due to the marginal benefits of interest tax shields. As the debt-to-value ratio further increases the marginal agency costs rise and the value of the firm begins to fall.

Consider first the possibility of lengthening the time to maturity on the bond while keeping the present value of the bond constant. In Panel A of Table VII we show this comparison. At an initial commodity price of \$0.45/pound and a firm with an outstanding bond with annual debt service payments totalling \$0.5 million for 5 years faces a high probability of bankruptcy. The levered firm has a value of \$1.602 million—compared to a first best value of \$6.9 million. The market value of the bond is \$1.374 million, a discount of 42% from par value. It is possible to design a bond with longer maturity, lower total debt service, and with an approximately equivalent market value: 15 years, \$0.047 annually, and current market value \$1.390. If this longer maturity bond is substituted or exchanged for the shorter maturity bond, then the value of the firm rises by over \$5 million or more than three times. While the tax shields on the two bonds are almost equal, the longer maturity bond has significantly lower agency costs, the difference being the difference in firm value of \$5 million. The shorter

## Table VII
## The Effect of Maturity Length on Agency Cost

| Panel A: Initial Commodity Price, $s = \$0.45/\text{pound}$ | | | | | |
|---|---|---|---|---|---|
| Final Maturity Date | Annual Bond Payments | Levered Firm Value $v^P(s, Q, t; 1)$ | Bond Value $b^P(s, Q, t; 1)$ | Tax Shields $\eta^P(s, Q, t; 1)$ | Agency Costs $\psi^P(s, Q, t; 1)$ |
| 5 years | 0.5 | 1.602 | 1.374 | 0.30 | 5.598 |
| 15 years | 0.047 | 6.836 | 1.390 | 0.31 | 0.374 |

| Panel B: Initial Commodity Price, $s = \$0.65/\text{pound}$ | | | | | |
|---|---|---|---|---|---|
| Final Maturity Date | Annual Bond Payments | Levered Firm $v^P(s, Q, t; 1)$ | Bond Value $b^P(s, Q, t; 1)$ | Tax Shields $\eta^P(s, Q, t; 1)$ | Agency Costs $\psi^P(s, Q, t; 1)$ |
| 5 years | 1.75 | 15.209 | 8.212 | 0.164 | 1.095 |
| 15 years | 0.95 | 13.581 | 8.207 | 0.478 | 3.037 |

maturity bond with its attendant high nominal debt service requirements puts the firm in imminent danger of bankruptcy. The equity holders' call option on the firm is way out of the money, and they are unlikely to want to continue to pay the debt service and maintenance costs in order to maintain their option. Only a very quick rise in prices will keep the firm out of costly bankruptcy. The longer maturity bond has lower debt service and an effectively longer time to maturity on the equity holders' call option. There is an attendant greater probability that their option will finish in the money, and consequently a lower probability of costly bankruptcy.

Extending the maturity of the bond and lowering the debt service payments does not always increase the value of the firm. When the firm is far from bankruptcy the owners will pay off the short-term bond and quickly return to a first best operating policy. In the meantime, if the term of the bond is short, then the total variance in price is relatively small, and the firm is likely to continue open, never making a suboptimal closing decision. If the maturity of the bond is lengthened, the firm will operate levered for a longer period of time, and the total variance in the price during this term is correspondingly greater. The firm is more likely to hit the close and open trigger prices repeatedly. While the firm is levered it will be opening and closing according to less than first best policy, lowering the value of the firm. The increased time to maturity also increases the possibility that the price will fall far enough to induce costly bankruptcy. A shorter maturity bond forces the equity holders to make the exercise decision relatively quickly when their option is in the money, and hence avoids this costly bankruptcy choice. This case is demonstrated in Panel B of Table VII. At an initial commodity price of $0.65/pound, and with an outstanding bond-paying annual debt service of $0.95 million for 15 years, the firm has a value $13.581

million. The bond has a market value of $8.207 million. The maturity of the bond could be shortened to 5 years with annual debt service payments totalling $1.75 million, and its value would be approximately the same at $8.212. The value of the firm would rise to $15.209 million because the agency costs would have fallen by nearly $2 million.

## IV. Conclusion

In this paper we have shown how to adapt a contingent claims model of the firm to reflect the incentive effects of the capital structure and thereby to measure the agency costs of debt. An important feature of our model is the existence of an underlying analysis of the firm and the stochastic features of its product market. We solve directly for the operating policy that is optimal for the equity owners and compare this with the first best operating policy. Our measure of the agency costs of debt is directly related to this underlying change in the use of the productive assets. The model determines the value of the firm and its associated liabilities incorporating the agency consequences of debt as well as the tax benefits.

Extending the agency literature to incorporate contingent claims techniques enables us to make a more refined use of the insights developed in that literature. The contingent claims technique yields a measure of agency costs that is robust to variations in the underlying parameters including the stochastic variable determining the firm's value. This measure can then be used to compare different capital structures and to analyze the agency effects under different circumstances facing the firm. We have already illustrated in this paper how the model can be used to compare the size of the agency costs associated with alternative maturity lengths of comparable debt instruments, depending upon the price of the firm's product. We also believe that this model can be extended to analyze debt contracts of fundamentally different design.

### REFERENCES

Black, F., and M. Scholes, 1973, The pricing of options and corporate liabilities, *Journal of Political Economy* 81, 637–654.
Brennan, M., and E. Schwartz, 1984, Optimal financial policy and firm valuation, *Journal of Finance* 39, 593–607.
————, 1986, Evaluating natural resource investments, *Journal of Business* 58, 135–157.
Fischer, E., R. Heinkel, and J. Zechner, 1989, Dynamic recapitalization policies and the role of call premia and issue discounts, *Journal of Financial and Quantitative Analysis* 24, 427–446.
Mello, A., and J. Parsons, 1991, Indexation and the optimal design of securities, Working paper, Massachusetts Institute of Technology.
Merton, R., 1974, On the pricing of corporate debt: The risk structure of interest rates, *Journal of Finance* 29, 449–470.
Mikkelson, W., and M. Partch, 1986, Valuation effects of security offerings and the issuance process, *Journal of Financial Economics* 15, 31–60.

Myers, S., 1977, Determinants of corporate borrowing, *Journal of Financial Economics* 5, 146–175.

Ross, S., 1978, A simple approach to the valuation of risky streams, *Journal of Business* 51, 453–475.

Shapiro, A., and S. Titman, 1986, An integrated approach to corporate risk management, in J. Stern and D. Chew, eds.: *The Revolution in Corporate Finance* (Basil Blackwell, New York).

# Daftar Pustaka

Günter Bamberg, Klaus Spremann (auth.) (1987), Agency Theory, Information, and Incentives, Springer Berlin Heidelberg

Barry Barnes (1999), Understanding Agency: Social Theory and Responsible Action, SAGE

Roger Frie (2008), Psychological Agency: Theory, Practice, and Culture,  A Bradford Book

Craig W. Gruber, Matthew G. Clark, Sven Hroar Klempe, Jaan Valsiner (eds.) (2015), Constraints of Agency: Explorations of Theory in Everyday Life, Springer International Publishing

Alexander Pepper (2019), Agency Theory and Executive Pay: The Remuneration Committee's Dilemma, Springer International Publishing,Palgrave Pivot

de Matos, Joao Amaro (2001); Theoritical Foundations of Corporate Finance; Princento University Press

Hart, Oliver (1995); Firms Contracts and Financial Structure; Oxford University Press.

Martin, John D., Cox, Samuel H. and R. D. Macminn (1988); The Theory of Finance: Evidence and Applications; The Dryden Press. (MCM)

Smith, C. W. (1990); The Modern Theory of Corporate Finance; McGraw Hill (CS)

Stern, J. M. and D. H. Chew (1986); The Revolution in Corporate Finance; Basil Blackwell (SC)

Dathine, J-P and J. B. Donalson (2014), Intermediate Fianancial Theory, Elsevier Academic Press (DD)

Tirole, J. (2006), The Theory of Corporate Finance, Princenton University Press (JT)

Adler Haymans Manurung (2012), Teori Keuangan perusahaan, PT Adler Manurung Press (AHM)

Arnold, Glen (2013), Corporate Financial Management, 5th Eds., PEARSON.

Copeland, T. E.,  Weston, J. F., & Shastri, K. (2013). Financial Theory and Corporate Policy: Pearson New International Edition. 4th Edition. Upper Saddle River, NJ: Pearson Education. ISBN: 1292034815, 9781292034812. (CWS)

Eichberger, Jurgen and Ian R. Harper (1997); Financial Economics; Orford University Press

Jones, Chris (2008); Financial Economics; Routledge.

Kettell, Brian (2001); Financial Economics: Making Sense of Market Information; Prentice Hall.

Brealey, R. A.; S. C. Myers and F. Allen (2014); Principles of Corporate Finance; McGraw Hill

Emery, Douglas R. and John D. Finnerty (1997), Corporate Financial Management; Prentice Hall.

Weston, J. F. and T. E Copeland (1986); Managerial Finance; The Dryden Press

Hawawini, G. and C. Viallet (2015), Finance for Executives: Managing for Value Creation; CENGAGE Learning

Manurung, Adler H (2021), Keuangan Perusahaan, PT Adler Manurung Press

Ruth Bender and Keith Ward (2002). Corporate Financial Strategy, 2nd eds., Butterworth Heinemann

Justin Pettit (2007), Strategic Corporate Finance: Application in Valuation & Capital Structure, John Wiley & Sons.

Rajesh Kumar (2017), Strategic Financial Management: Case Book; Academic Press

2011 H. Kent Baker, Gerald S. Martin(auth.) - Capital Structure and Corporate Financing Decisions_ Theory, Evidence, and Practice

J. Fred Weston and Thomas E. Copeland (1986), Managerial Finance, 8th Eds. The Dryden Press

Dathine, J-P and J. B. Donalson (2014), Intermediate Fianancial Theory, Elsevier Academic Press (DD)

Emery, Douglas R. and John D. Finnerty (1997), Corporate Financial Management; Prentice Hall

Haley, C. W. and L. D. Schall (1979); The Theory of Financial Decisions; McGraw Hill.

Meggison, W. L. (1997), Corporate Finance Theory; Addison Wesley

J. F. Weston and T. E Copeland (1986); Managerial Finance; The Dryden Press

Brealey, R. A.; S. C. Myers and F. Allen (2014); Principles of Corporate Finance; McGraw Hill

# Riwayat Hidup Penulis

**Adler Haymans Manurung**, dilahirkan di Porsea, Tapanuli Utara pada 17 Desember tahun 1961. Pendidikan Sekolah Dasar (SD) sampai Sekolah Menengah Atas di Medan. Selanjutnya, pendidikan perguruan tingginya dimulai dari Akademi Ilmu Statistik dengan lulus Ranking Pertama pada tahun 1983. Sarjana Ekonomi (SE) diperolehnya dari Program Extension Fakultas Ekonomi Universitas Indonesia pada tahun 1987. Pendidikan program S2 dengan gelar Master of Commerce (M.Com) dari University of Newcastle, Australia pada tahun 1995 dan Magister Ekonomi (ME) dari Fakultas Ekonomi Universitas Indonesia pada tahun 1996. Doktor dalam bidang Keuangan diperoleh dari FEUI pada 17 Oktober 2002 dengan predikat "Cum-Laude". Lulus Sarjana Hukum dengan menekuni Hukum Ekonomi dari Fakultas Hukum Universitas Kristen Indonesia pada tahun 2007. Adler juga telah menyelesaikan Kursus Pajak Brevet A dan B di STAN, Jakarta pada tahun 2007.

Dalam Bidang Bisnis, Adler saat ini mengelola beberapa perusahaan, President Direktur PT Valuasi Investindo, PT Finansial Bisnis Informasi, dan PT Adler Manurung Press. Juga menjadi Komisaris PT Rygrac Capital dan PT Putra Nauli (bergerak dalam bidang pupuk kompos di Porsea – Kabupaten Tobasa, SUMUT) dan Ketua Dewan Pembina Yayasan Tobasa Membangun. Sebelumnya, Adler bergabung dengan PT Nikko Securities Indonesia pada periode Nopember 1996 sampai April 2010 dengan jabatan Direktur Fund Management dan dimana sebelumnya bekerja pada PT BII Lend Lease Investment Services sebagai Associate Direktur Riset sejak Maret 1995 sampai dengan Oktober 1996 dan sebagai Senior Manager Research Analyst pada Lend Lease Corporate Services, Australia, sejak Juli 1994. Sebagai Fund Manager telah mengalami asam garam dan saat ini telah mengelola dana diatas Rp. 2 trilliun. Investor yang sangat mengenalnya menyebut **pelindung dana investor** karena sangat hati-hatinya. Adler memulai karir dalam pasar modal pada tahun

1990 dan bekerja sebagai Research Analyst di perusahaan sekuritas. Pada periode 2010 – 2014 menjadi Ketua Komite Tetap Fiskal dan Moneter, Kadin Indonesia.   Adler telah menulis buku sebagai berikut:

1. Statistik Lanjutan (Advanced Statistics Problem) Penerbit : Universitas Tarumanegara (1989).
2. Teknik Peramalan Bisnis dan Ekonomi (Forecasting Method for Business and Economic) Penerbit: PT. Rineka Cipta (1990)
3. Pengambilan Keputusan; Pendekatan Kuantitatif (Decision Theory; Quantitative and Economic) Penerbit: PT. Rineka Cipta (1991)
4. Analisis Saham Indonesia (Stock Analysis in Indonesia) Penerbit: Economic Student's Group (1992)
5. Lima Bintang untuk Agen Penjual Reksa Dana, Penerbit: Ghalia Indonesia, 2002.
6. Memahami Seluk Beluk Instrumen Investasi. Penerbit: PT Adler Manurung Press, April - 2003
7. Berinvestasi, Pendirian dan Pembubaran Reksa Dana: Pegangan untuk Manajer Investasi dan Investor; Penerbit: PT Adler Manurung Press, Agustus – 2003.
8. Pasar Keuangan & Lembaga Keuangan Bank & Bukan Bank; Penerbit: PT Adler Manurung Press, Agustus 2003. (Sebagai Penulis Ketiga)
9. Strategi Memenangkan Transaksi Saham di Bursa (Strategic to win stock transaction in Bourse), PT Elex Media Komputindo (Gramedia Group); Agustus 2004.
10. Penilaian Perusahaan (Company Valuation); Penerbit: PT Adler Manurung Press, September 2004 – diperbaharui dengan Judul "Valuasi Wajar Perusahaan".
11. Dasar-dasar Keuangan Bisnis: Teori dan Aplikasi; Penerbit: PT Elex Media Komputindo, Jakarta, Mei 2005., (Penulis Kedua dari tiga Penulis)
12. Wirausaha: Bisnis UKM, Kompas Agustus 2005
13. Ke Arah Manakah Bursa Indonesia dibawa?, Penerbit: PT Elex Media Komputindo, Jakarta Oktober 2005
14. Ekonometrika: Teori dan Aplikasi; PT Elex Media Komputindo, Jakarta Desember 2005. (Penulis Kedua dari tiga penulis)
15. Ke Mana Investasi ? Kiat dan Panduan Investasi Keuangan di Indonesia; Penerbit Buku Kompas, Maret 2006.

16. Dasar-Dasar Investasi Obligasi; PT Elex Media Komputindo; Mei 2006.
17. Aktiva Derivatif: Pasar Uang, Pasar Modal, Pasar Komoditi, dan Indeks; PT Elex Media Komputindo; Desember 2006, (Penulis Kedua)
18. Cara Menilai Perusahaan; PT Elex Media Komputindo; Januari 2007,
19. Sekuritisasi Aset, PT Elex Media Komputindo, Maret 2007
20. Wanita Berbisnis UKM – Makanan, Kompas Maret 2007
21. Pengelolaan Portofolio Obligasi, PT Elex Media Komputindo, April 2007
22. Reksa Dana Investasiku, Kompas September 2007.
23. Pendanaan UKM, Kompas Januari 2008.
24. Financial Planner, Kompas, Maret 2008
25. Obligasi: Harga, dan Perdagangannya, ABFI Institute Perbanas, Januari 2009. Direvisi dan diterbitkan PT Adler Manurung Press, 2011.
26. Ekonomi Keuangan dan Kebijakan Moneter; Penerbit Salemba Empat, 2009 (Penulis Kedua, dengan Dr. Jonni Manurung)
27. Successful Financial Planner: A Complete Guide, PT Gramedia Widiasarana Indonesia, Agustus 2009
28. Kaya dari Bermain Saham; Penerbit Buku Kompas, Oktober 2009 (Di Revisi pada Maret 2021).
29. Metode Riset: Keuangan dan Investasi Empiris, ABFI Institute Perbanas Press, November 2009 – Bersama Wilson R. L. Tobing Ph.D.
30. Sukses Menjual Reksa Dana, PT Grasindo, 2010
31. Kaya dari Bermain Opsi; Penerbit Buku Kompas, 2010
32. Ekonomi Finansial; PT Adler Manurung Press, Jakarta, 2010
33. Metode Penelitian: Keuangan, Investasi dan Akuntansi Empiris; PT Adler Manurung Press, Mei 2011, diperbaiki dan diterbitkan Kembali pada tahun 2019 dengan penulis kedua Dr. Dyah Budiastuti.
34. Restrukturisasi Perusahaan: Merger, Konsolidasi, Merger dan Akuisisi serta Pembiayaannya, PT Adler Manurung Press, Agustus 2011
35. Teori Keuangan Perusahaan; PT Adler Manurung Press, Januari 2012
36. Teori Investasi: Konsep dan Empiris; PT Adler Manurung Press, Agustus 2012.

37. Investasi dan Manajemen Portofolio, Modul untuk FE Universitas Terbuka, 2012
38. Initial Public Offering (IPO): Konsep, Teori dan Proses; PT Adler Manurung Press, April 2013
39. Otorias Jasa Keuangan: Pelindung Investor; PT Adler Manurung Press, September 2013.
40. Berani Bermain Saham, Buku Kompas, September 2013.
41. Pasar Futures Indonesia: Tradisional to Finansial; PT Adler Manurung Press, Agustus 2014.
42. Pengukuran Risiko, PT Adler Manurung Press, Oktober 2014
43. Manajemen Treasuri: Dasar dan Instrumen; PT Adler Manurung Press, 2015
44. Konstruksi Portofolio Efek di Indonesia; PT Adler Manurung Press, Februari 2016
45. Raja Manurung tu Tuan Sogar Manurung dan Pomparannnya: "Mulak Ma Ogung tu Sakke Na; Jakarta: PT Adler Manurung Press, September 2016
46. Cadangan Devisa dan Kurs Valuta Asing; Buku Kompas, Oktober 2016
47. Manajemen Risiko Finansial: Perbankan, PT Adler Manurung Press, Februari 2017. Telah direvisi dengan judul "Manajemen Risiko Finansial untuk Industri Jasa Keuangan" ditulis Mohammad Hamsal, Adler Haymans Manurung, Benny Hutahayan dan Jenry Cardo Manurung.
48. Manajemen Aset dan Liabilitas, PT Adler Manurung Press, Juni 2017
49. Model dan Estimasi dalam Riset Manajemen dan Keuangan; PT Adler Manurung Press, Juli 2019.
50. Enterprise Risk Management, PT Adler Manurung Press, Jakarta, Februari 2020.
51. Bank Business Performance, PT Adler Manurung Press, Nopember 2020, Penulis Pertama dari 4 Penulis (Benny Hutahayan, Kevin Deniswara dan Tipri Rose Kartika)
52. Investasi: Teori dan Empiris; PT Adler Manurung Press, Nopember 2020
53. Manajemen: Teori dan Perkembangannya, PT Adler Manurung Press, Februari 2021
54. Keuangan Perusahaan, PT Adler Manurung Press, Juli 2021
55. Financial Modeling: Microsoft Excel, PT Adler Manurung Press, Februari 2022.

56. Regression and Extension, PT Adler Manurung Press, Maret 2022
57. Market Microstructure: A Reading, PT. Adler Manurung Press, May 2022

Disamping sebagai penulis buku, Adler juga aktif sebagai kolumnis dalam bidang pasar modal diberbagai surat kabar, majalah nasional serta majalah internasional serta **pengasuh kolom Investasi di Harian Kompas Minggu**. Tulisan penelitian empirisnya dapat dibaca pada Jurnal terkemuka di Indonesia, seperti Jurnal Riset dan Akuntansi Indonesia (JRAI), Jurnal Kelola dari UGM dan Management Usahawan dari FEUI serta Jurnal Perbankan dari STIE Perbanas. Disamping itu, Adler juga menjadi pembicara dalam konferensi ilmiah internasional dan juga menjadi staf pengajar pada MM-FEUI, Pascasarjana FEUI; Doktor Bisnis di MB – IPB dan Program Doktor Manajemen Bisnis, Universitas Padjadjaran, Bandung dan Pascasarjana ABFI Institute Perbanas; Magister Manajemen – Universitas Negeri Jakarta serta Fakultas Ekonomi – Universitas Tarumanagara. Kepangkatan penulis dalam mengajar dari Departemen Pendidikan yaitu "**Professor**" pada tahun 2008 dalam bidang Investasi, Pasar Modal, Keuangan dan Perbankan dengan dengan Surat Keputusan Menteri Pendidikan Nasional Republik Indonesia Nomor: 77548/A4.5/KP/2008, tertanggal 1 Desember 2008. Adler telah ditugaskan BAN-PT sebagai Assessor BAN-PT. Penulis juga menjadi Chief Editor Journal Keuangan dan Perbankan yang diterbitkan ABFI Institute Perbanas dan merupakan satu dari lima jurnal terakreditasi B di Dirjen Perguruan Tinggi. Adler telah memperoleh ijin sebagai Wakil Manajer Investasi dan Wakil Penjamin Emisi Efek dari Bapepam. Penulis juga memperoleh gelar professional Chartered Financial Consultant (ChFC) dan Chartered Life Underwriting (CLU) dari American College serta Registered Financial Consultant (RFC) dari International Association of Registered Financial Consultant, Agustus 2004. Adler juga memiliki sertifikasi Eksekutif Risk Management Corporate Professional (ERMCP) pada tahun 2009 dari ERMI - Singapore. Penulis juga aktif dalam bidang organisasi sebagai Ketua Assosiasi Pengelola Reksa Dana Indonesia (APRDI) pada periode 2001 – 2004. Saat in penulis menjadi Technical Advisor pada Internasional Association of Registered Financial Consultant for Indonesia. Pada tahun 2004, penulis masuk nominasi 10 besar "The Most Popular Analyst "dan memperoleh "The Most Popular Analyst 2005" atas survey **Frontier**

**Indonesia.** Adler juga menjadi salah satu juri di REBI (Recognize Bisnis) yang dikoordinir Koran Sindo dan Frontier.

Sejak September 2012, Prof. Adler H. Manurung diangkat menjadi Guru Besar Pasar Modal, Investasi, Keuangan dan Perbankan pada Sampoerna School of Business (SSB) dan kemudian 1 September 2012 menjadi Kepala Program Studi Manajemen dan sejak 1 Mei 2013 diangkat Putera Sampoerna Foundation menjadi Ketua STIE Putera Sampoerna dan kemudian menjadi Dekan Fakultas Bisnis, Universitas Siswa Bangsa Internasional (USBI). Jurnal Bisnis dan Kewirasusahaan dibangun di SSB dan sudah terbit dan beredar bagi para akademisi maupun praktisi. Jabatan Ketua STIE Putera Sampoerna berakhir pada 30 April 2014. Menjadi adviser PT Bursa Berjangka Jakarta sejak 1 Juli 2013 sampai sekarang dalam rangka membuat produk Bonds *Futures*. Prof. Dr. Adler H. Manurung diangkat menjadi Dosen Tetap dan sekaligus Guru Besar Pasar Modal, Investasi dan Perbankan di Fakultas Ekonomi Universitas Bina Nusantara, Jakarta sejak 1 Nopember 2014. Februari 2021 menjadi Guru Besar Pasar Modal dan Perbankan Universitas Bhayangkara Jakarta Raya dan mendirikan Program Studi Doktor Ilmu Manajemen. Sejak Oktober Tahun 2013 mendirikan Assosiasi Analis Pasar Investasi dan Perbankan dan menjadi Presiden assosiasi ini, dimana assosiasi ini memberikan sertfikasi professional dengan gelar CIMBA. Penulis juga telah menyelesaikan Pendidikan Kepemimpinan Nasional, PPSA-XX, Lemhanas 2015. Sejak 2016, mulai mengajar di Universitas Pertahanan (UNHAN) dibawah Kementerian Pertahanan (KEMENHAN).

Prof. Dr. Adler Haymans Manurung menikah dengan Ir. Marsaurina Yudiciana boru Sitanggang pada tahun 1990. Atas pernikahan tersebut dikaruniai anak dua orang yaitu Castelia Romauli dan Adry Gracio. Castelia Romauli sudah menyelesaikan kuliah di Universitas Negeri Jakarta dan sedang mengikuti kuliah Pascasarjana di Atmajaya dan bekerja pada Bank Internasional. Adry Gracio telah lulus dari Jurusan Ilmu Ekonomi di FEUI dengan predikat Cum-Laude, serta juga telah lulus Master of Science dari London School Economics – UK dan saat ini sudah bekerja.

**Jhonni Sinaga,** kelahiran Padang Sidempuan, Sumatera Utara, pada 20 Desember 1968. Putra dari ayah dan ibu yang berprofesi sebagai guru ini menamatkan pendidikan dasar hingga menengah atas di Kotamadya Sibolga, Sumatera Utara.

Pada tahun 1990, menyelesaikan pendidikan Bahasa Inggris di Universitas Sumatera Utara. Studi Ekonomi Manajemen juga ditempuh dan diselesaikannya di universitas yang sama pada tahun 1994. Pada tahun 2014, ia menyelesaikan Program Magister Manajemen (S-2; M.M.) di Universitas Mulawarman. Studi Program Doktor Manajemen (S-3) juga ditempuh dan diselesaikannya di universitas yang sama pada tahun 2019.

Pada tahun 1996 – Maret 2000, ia bekerja sebagai Kepala Seksi Akuntansi dan Keuangan pada Salim Plantations (Indofood Plantations, Tbk.). Semenjak April 2000, ia dipercaya sebagai *Head of Internal Audit Department* pada B.W. Plantations, Tbk. Pada perusahaan ini karirnya meningkat pada Juli 2003, ia diangkat menjadi *Accounting and Tax Manager*. Per Desember 2005, ia menjabat *Head of Internal Audit Department* pada REA Kaltim Plantations Group (*Subsidiary of REA Holding, a U.K. Public Listed Company at London Stock Exchange*). Sebelas tahun kemudian, mulai Agustus 2016, ia menjadi Head of Operation Finance and Accounting. Pada 01 Agustus 2019, ia resmi mengundurkan diri dari REA Kaltim Plantations Group. Pada April 2020 mendirikan lembaga yang bergerak dalam bidang konsultasi manajemen perusahaan dengan bendera J. J. Manajemen Konsultasi dan pada Nopember 2020 resmi mendirikan entitas yang bergerak dalam bidang transportasi (*trucking and logistic*) dengan bendera PT. JeJe Harapan Transindo (JeJe Trans Group). Kedua bidang usaha ini tumbuh dan berkembang melampaui ekspektasi hingga saat ini.

Semenjak Agustus 2021 menjadi dosen tidak tetap di Universitas Kristen Indonesia (UKI) dan Januari 2022 diangkat sebagai dosen tetap di Universitas Bhayangkara Jakarta Raya hingga saat ini. Selama bekerja sebagai profesional perusahaan, ia aktif sebagai tutor untuk materi-materi seperti *Internal Control*, *Good Corporate*

*Governance* (GCG), *Supervision Management*, *Coaching for Performance*, *Motivation*, *Budgetting*, dan *Finance and Accounting*. Pada tahun 2006, ia berhasil menciptakan konsep usaha perkebunan kelapa sawit baru dengan nama "**PRO EXISTENCE**".