# H I E R

# Harvard Institute of Economic Research

When Is Reputation Bad?

by

*Jeffrey Ely, Drew Fudenberg
and David K. Levine*

*June 2002*

Harvard University
Cambridge, Massachusetts

# When is Reputation Bad?[1]

Jeffrey Ely

Drew Fudenberg

David K. Levine[2]

First Version: April 22, 2002

This Version: June 12, 2002

In traditional reputation theory, reputation is good for the long-run player. In "Bad Reputation," Ely and Valimaki give an example in which reputation is unambiguously bad. This paper characterizes a more general class of games in which that insight holds, and presents some examples to illustrate when the bad reputation effect does and does not play a role. The key properties are that participation is optional for the short-run players, and that every action of the long-run player that makes the short-run players want to participate has a chance of being interpreted as a signal that the long-run player is "bad. We also broaden the set of commitment types, allowing many types, including the "Stackelberg type" used to prove positive results on reputation. Although reputation need not be bad if the probability of the Stackelberg type is too high, the relative probability of the Stackelberg type can be high when all commitment types are unlikely.

## 1. Introduction

A long-run player playing against a sequence of short-lived opponents can build a reputation for playing in a specific way and so obtain the benefits of commitment power. To model these "reputation effects," the literature following Kreps and Wilson [1982] and Milgrom and Roberts [1982] has supposed that there is positive prior probability that the long-run player is a "commitment type" who always plays a specific strategy.[3] In "Bad Reputation," Ely and Valimaki [2001] (henceforth EV) construct a striking example in which introducing a particular commitment type hurts the long-run player. When the game is played only once and there are no commitment types, the unique sequential equilibrium is good for the long-run player. This remains an equilibrium when the game is repeated without commitment types, regardless of the player's discount factor. However, when a particular "bad" commitment type is introduced, the only Nash equilibria are "bad" for a patient long-run player.[4]

What is not clear from EV is when reputation is bad. This paper extends the ideas in EV to a more general class of games in an effort to find the demarcation between "bad" and "good" reputation. In addition, we try to relate the EV conclusions to past work on reputation effects.

Reputation effects are most powerful when the long-run player is very patient, and Fudenberg and Levine [1992] (FL) provided upper and lower bounds on the limiting values of the equilibrium payoff of the long-run player as that player's discount factor tends to 1. The upper bound

---

[3] See Sorin [1999] for a recent survey of the reputation effects literature, and its relationship to the literature on merging of opinions.

[4] It is obvious that incomplete information about the long-run player's type can be harmful when the long-run player is impatient, since incomplete information can be harmful in one-shot games. Fudenberg-Kreps [1987] argue that a better measure of the "power of reputation effects" is to hold fixed the prior distribution over the reputation-builder's types, and compare the reputation-building scenario to one in which the reputation builder's opponents do not observe how the reputation builder has played against other opponents. They discuss why reputation effects might be detrimental in the somewhat different setting of a large long-run player facing many simultaneous small but long run opponents.

corresponds to the usual notion of the "Stackelberg payoff." The lower bound, called the "generalized Stackelberg payoff," weakens this notion to allow the short-run players to have incorrect beliefs about the long-run player's strategy, so long as the beliefs are not disconfirmed by the information that the short-run players get to observe. When the stage game is a one-shot simultaneous-move game, actions are observed, and payoffs are generic, these two bounds coincide, so that the limit of the Nash equilibrium payoffs as the long-run player's discount factor tends to one is the single point corresponding to the Stackelberg payoff. For extensive-move stage games, with public outcomes corresponding to terminal nodes, the bounds can differ. However, although FL provided examples in which the lower bound is attained, in those examples the upper bound was attained as well, and we are not aware of past work that determines the range of possible limiting values for a fairly general class of games.

Here we examine the upper bound more closely for a specific class of games designed to capture the insight of EV. Specifically, we define a class of "bad reputation" games, in which the long-run player can do no better than if the short-run players choose not to participate. This extends the EV example in a number of ways. We allow a broad class of stage games in which participation by the short-run players is optional; allowing for many actions, many signals, many short-run players, and a wide variety of payoffs. Especially important, we allow for a broad range of types, including types that are committed to "good" actions, as well as types that are committed to "bad" actions. Earlier research suggests that to attain the upper bound on the long-run player's payoff, it can be important to include the "Stackelberg type" that is committed to the stage-game action the long-run player would choose in a Stackelberg equilibrium.[5] We

---

[5]  EV consider two specifications for the bad type, either "committed" (to playing the bad action) or "strategic" (willing to play a different action occasionally to increase entry and the future payoff from playing "bad.") In a related model, Mailath and Samuelson [1998] argue that "bad" types – and  specifically strategic bad types – are more plausible than Stackelberg types. We are sympathetic to the argument that strategic bad types may be

find that the EV results fails if the probability of this Stackelberg type is too high, but extends to the case where the probability of the Stackelberg type is sufficiently low, but nonzero. This shows that it is not essential to rule out the types that support "good" reputation effects in order to derive the bad reputation result.

By extending the EV example to a broad class of stage games we are able to more clearly identify the types of assumptions key to a bad reputation. There are several such properties, notably that the short-run players can either individually or collectively choose not to participate. However, most of the assumptions on the structure of the game seem to involve little loss of economic applicability. The key substantive assumption seems to be that every action of the long-run player that makes the short-run players want to participate in the game has a chance of being misconstrued as a signal of a "bad reputation."

EV motivate their example by considering an automobile mechanic who has specialized knowledge of the work that needs to be done to repair the car. We think that we have identified a broader class of bad reputation games that can be interpreted as "expert advice." This includes consulting a doctor or stockbroker, or in the macroeconomics context, can be the decision whether or not to turn to the IMF for assistance. In EV, the short-run players observe only the advice, but not the consequences of the advice. Here we explicitly consider what happens when the short-run players observe the consequences as well. We also show that there are other distinct classes of games with rather different observation structures that are bad reputation games, such as our "teaching evaluation" game, where "advice" is not an issue because the long-run player does not privately observe anything that is payoff-relevant for the short-run player. Finally, we illustrate the boundaries of bad reputation by giving a number

---

more likely than commitment types, but this does not imply that the probability of commitment types should be zero. Instead, we would argue that it is preferable for models to allow for a wide range of types, especially those with fairly simple behavior rules.

of examples and classes of participation games that are not bad reputation games.

## 2. The Model

### *2.1. The Dynamic Game*

There are $N + 1$ players, a long run-player 1, and $N$ short-run players $2...N + 1$. The game begins at time $t = 1$ and is infinitely repeated. Each period, each player $i$ chooses from a finite action space $A^i$. We denote individual actions $a^i$, and action profiles by $a$. We also use $a^{-i}$ to denote the play of all players except player $i$ and $a^{-i-j}$ to denote the play of all players except players $i$ and $j$.

The long-run player discounts the future with discount factor $\delta$. Each short-run player plays only in one period, and is replaced by an identical short-run player in the next period. There is a set $\Theta$ of types of long-run player. There are two sorts of types: type $0 \in \Theta$ is called the "rational type," and is the focus of our interest, with utility described below. For each pure action $a^1$, type $\theta(a^1)$ is a "committed type," that is constrained to play $a^1$. These are the only possible types in $\Theta$. Note that we do not require that every pure action commitment type has positive probability. The stage game utility functions are $u^i(a)$, where $u^1(a)$ corresponds to the long-run player of type $\theta = 0$. The common prior distribution over long-run player types is denoted $\mu(0)$.

There is a finite public signal space $Y$ with signal probabilities $\rho(y \mid a)$. All players observe the history of the public signals. Short-run players observe only the history of the public signals, and in particular observe neither the past actions of the long-run player, nor of previous short-run players. We do not assume that the payoffs depend on the actions only through the signals, so the short-run players at date $t$ need not

know the realized payoffs of the previous generations of short-run players.[6]

We let $h_t = (y_1, y_2 \ldots, y_t)$ denote the public history through the end of period $t$. We denote the null history by $0$. We let $h_t^1$ denote the private history known only to the long-run player. This includes his own actions, and may or may not include the actions of the short-run players he has faced in the past. A strategy for the long-run player is a sequence of maps $\sigma^1(h_t, h_t^1, \theta) \in \text{conhull}\, A^1 \equiv \mathcal{A}^1$; a strategy profile for the short-run players is a sequence of maps $\sigma^{-1}(h_t) \in \times_{j \neq 1} \text{conhull}\, A^j \equiv \mathrm{A}^{-1}$. A short-run profile $\alpha^{-1}$ is a Nash response to $\alpha^1$ if $u^i(\alpha^1, \alpha^i, \alpha^{-1-i}) \geq u^i(\alpha^1, a^i, \alpha^{-1-i})$ for all $a^i \in A^i$. We denote the set of short-run Nash responses to $\alpha^1$ by $B(\alpha^1)$.

Given strategy profiles $\sigma$, the prior distribution over types $\mu(0)$ and a public history $h_t$ that has positive probability under $\sigma$, we can calculate from $\sigma^1$ the conditional probability of long-run player actions $\bar{\alpha}^1(h_t)$ given the public history. A *Nash Equilibrium* consists of strategy profiles $\sigma$ such that for each positive probability history

1) $\sigma^{-1}(h_t) \in B(\bar{\alpha}^1(h_t))$ [short-run players optimize]
2) $\sigma^1(h_t, h_t^1, \theta(a^1)) = a^1$ [committed types play accordingly]
3) $\sigma^1(\cdot, \cdot, 0)$ is a best-response to $\sigma^{-1}$ [rational type optimizes].

Given a Nash equilibrium, and a positive probability history $h_t$ let $v^1(h_t)$ denote the expected continuation value to the rational long-run player, and let $\mu(h_t), \alpha^{-1}(h_t)$ be the posterior beliefs and strategy of the short-run players conditional on the history. Notice that the expected average present value to the rational long-run player conditional on a positive probability public/private history pair must not depend on the private history $h_t^1$, or the rational long-run type would be failing to optimize. If $a^1$ has positive probability under $\bar{\alpha}^1(h_t)$, and $a^{-1}$ positive probability under $\alpha^{-1}(h_t)$, then we define

---

$$v^1(h_t, a) \equiv (1 - \delta)u^1(a) + \delta\sum_y \rho(y \mid a)v^1(h_t, y) \ .$$

When $\alpha^1$ and $\alpha^{-1}$ put weight only on such positive-probability $a^1, a^{-1}$, it is convenient also to define $v^1(h_t, \alpha)$.

### 2.2  The Ely-Valimaki Example

We will use the EV example to illustrate our assumptions and definitions.  In EV, an action by the long-run player is a map from a privately observed signal $\omega \in \{E, T\}$ to announcements $\{e, t\}$; thus the long-run player's action space consists of $A^1 = \{ee, et, te, tt\}$. The two signals are i.i.d. and equally likely.  There is one short-run player each period who chooses an element of $A^2 = \{In, Out\}$.  The public signal takes on the values $Y = \{e, t, Out\}$. If the short-run player chooses $Out$ the signal is $Out$, that is $\rho(Out \mid a^1, Out) = 1$; otherwise the signal is the announcement of the long-run player, so $\rho(e \mid (et, In)) = \rho(e \mid (te, In)) = 1/2$, $\rho(e \mid (ee, In)) = 1$, and $\rho(e \mid (tt, In)) = 0$.  If the short-run player chooses $Out$, each player gets utility 0. If he plays $In$ and the long-run player's announcement is truthful (that is, matches the state), the short-run player receives $u$; if it is untruthful, it is $-w$ where $w > u > 0$. The "rational type" of long-run player has exactly the same stage-game payoff function as the short run players.  Thus when the long-run player is certain to be the rational type, the strategic form of the stage game is

|      | *In*                        | *Out* |
|------|-----------------------------|-------|
| *ee* | $(u - w)/2, (u - w)/2$       | $0, 0$ |
| *et* | $u, u$                      | $0, 0$ |
| *te* | $-w, -w$                    | $0, 0$ |
| *tt* | $(u - w)/2, (u - w)/2$       | $0, 0$ |

Figure 1

When the rational type is the only type in the model, there is an equilibrium where he chooses the action that matches the state, all short-run players enter, and the rational type's payoff is $u$. However, EV show that when there is also a probability that the long-run player is a "bad type" who always plays $ee$, the long-run player's payoff is bounded by an amount that converges to 0 as the discount factor goes to 1. The intuition for this result has three steps. First of all, the short-run players will not enter if the long-run player is too likely to play $ee$. Second, from Bayes rule it follows that there is some number $K$ such that $K$ successive observations of $E$ will make the posterior probability of the bad type so high that all subsequent short-run players play out. Third, when there have been $K - 1$ successive observations of $E$, the rational type of long run player is tempted to play *tt* instead of *et*, even though this lowers his short-run payoff, to avoid driving out the short-run players with another observation of *E*. Thus, the long-run player is tempted to take an action that is worse for both himself and the short-run players in order to avoid being incorrectly tagged as a "bad type." Our result will generalize this idea of a "temptation."

### 2.3. Participation Games and Bad Reputation Games

We consider "participation games" in which the short-run players may choose not to participate. The crucial aspect of non-participation by the short-run players is that it conceals the action taken by the long-run player from subsequent short-run players; this is what allows the lower bound on the long-run player's Nash equilibrium payoff in the EV example to be lower than Stackelberg payoff. We will then define "bad reputation" games as a subclass of participation games that have the additional features needed for the bad reputation result.

To model the option to not participate, we assume that certain public signals $y^e \in Y^e$ are *exit signals*. Associated with these exit signals are *exit profiles*, which are pure action profiles $e \in E^{-1} \subseteq A^{-1}$ for the short

run players. For each such $e$, $\rho(y^e \mid a^1, e) = \rho(y^e \mid e)$ for all $a^1$, and $\rho(Y^e \mid e) = 1$. In other words, if an exit profile is chosen, an exit signal must occur, and the distribution of exit signals is independent of long-run player action. Moreover, if $a^{-1} \notin E^{-1}$ then $\rho(y^e \mid a^1, a^{-1}) = 0$ for all $a^1 \in A^1, y^e \in Y^e$. We refer to $A^{-1} - E^{-1}$ as the entry profiles. Note that an entry profile cannot give rise to an exit signal. A *participation game* is a game in which $E^{-1} \neq \varnothing$. The remainder of the paper specializes to participation games.

We begin by distinguish actions by the long-run player that cause the short-run players to exit (unfriendly actions), and those that are needed to get them to enter (friendly actions).

**Definition 1:** *A finite non-empty set of pure actions $\widehat{A}^1$ for the long-run player is* unfriendly *if there is a number $\widehat{\alpha} < 1$ such that $\alpha^1(\widehat{A}^1) \geq \widehat{\alpha}$ implies $B(\alpha^1) \subseteq$ conhull $E^{-1}$.*

*Remark:* This definition says that unfriendly actions induce exit, in the strong sense that exit is the only best response if the probability of the unfriendly actions is sufficiently high. There will often be many sets of unfriendly actions. In the EV example the set $\{ee, tt, te\}$ is unfriendly, and so is any subset.

**Definition 2:** *A finite set of mixed actions $F^1$ for the long run player is* friendly *if there is a number $\underline{\alpha} > 0$ such that $B(\alpha^1) \cap \left[ A^{-1} - \text{conhull}(E^{-1}) \right] \neq \varnothing$ implies $\alpha^1 \geq \underline{\alpha} f^1$ for some $f^1 \in F^1$.*

*Remark:* This definition says that the probabilities given to every pure action must be bounded below by a scale factor times some friendly mixture if the short-run players are not to exit. Note that weight on a friendly action is necessary for entry, but need not be sufficient for entry. There may also be many different friendly sets. Suppose that $F^1$ is friendly of size $\widehat{\alpha}_0$, and let $\underline{\alpha} = \min\{f^1(a^1) > 0 \mid f^1 \in F^1, a^1 \in A^1\}$. Then if $f^1 \in F^1$ it may be replaced by any mixture over the support of $f^1$, and the resulting set will be friendly of size $\widehat{\alpha}_0 \underline{\alpha}$. Similarly, if we

have a friendly set and we eliminate mixtures $\tilde{f}^1 \in F^1$ whose support contained in the support of some different $f^1 \in F^1$, we get a new friendly set with a smaller value of $\underline{\alpha}$. In the EV example, the action $et$ is friendly, with

$$\underline{\alpha} = \frac{w - u}{w + u/2}.$$

Finally, consider a pure action $a^1$ such that $B(a^1) \cap \left[ A^{-1} - \text{conhull}(E^{-1}) \right] \neq \varnothing$. Since $a^1$ is pure, $a^1 \geq \underline{\alpha} f^1$ is possible only if $f^1 = a^1$. In other words, any pure action that permits short-run entry (such as $et$ in the EV example) *must* be in *every* friendly set. Moreover, if there is a single pure action that permits entry (again $et$) then this action is the unique friendly set, even if some mixed actions allow entry as well.

**Definition 3:** *The* support $A^1(F^1)$ *of a friendly set* $F^1$ *is the actions that are played with positive probability:* $A^1(F^1) \equiv \{a^1 \in A^1 \mid f^1(a^1) > 0, f^1 \in F^1\}$. *We say that a friendly set* $F^1$ *is* orthogonal *to an unfriendly set* $\widehat{A}^1$ *if* $\widehat{A}^1 \cap A^1(F^1) = \varnothing$.

Next we consider what signals may reveal about actions.

**Definition 4:** *We say that a set of signals* $\widehat{Y}$ *is* unambiguous *for a set of actions* $\widehat{A}^1$ *if for all* $a^{-1} \notin E^{-1}, \hat{y} \in \widehat{Y}, \hat{a}^1 \in \widehat{A}^1, a^1 \notin \widehat{A}^1$ *we have* $\rho(\hat{y} \mid \hat{a}^1, a^{-1}) > \rho(\hat{y} \mid a^1, a^{-1})$.

Notice that this is a strong condition: every action in $\widehat{A}^1$ must assign a higher probability to each signal in $\widehat{Y}$ than any action not in $\widehat{A}^1$. A given set of actions may not have signals that are unambiguous. In the case of the EV example, $E$ is an unambiguous signal for the sets $\{ee\}, \{ee, et, te\}$. The set $\{et\}$ does not have an unambiguous signal.

**Definition 5:** *An action* $a^1$ *is* vulnerable to temptation relative to a set of signals $\widehat{Y}$ *if there exists numbers* $\underline{\rho}, \tilde{\rho} > 0$ *and an action* $b^i$ *such that*

1)  If $a^{-1} \notin E^{-1}$, $\hat{y} \in \hat{Y}$, then $\rho(\hat{y} \mid b^1, a^{-1}) \leq \rho(\hat{y} \mid a^1, a^{-1}) - \underline{\rho}$.

2)  If $a^{-1} \notin E^{-1}$ and $y \notin \hat{Y} \cup Y^e$,
then $\rho(y \mid b^1, a^{-1}) \geq (1 + \tilde{\rho}) \rho(y \mid a^1, a^{-1})$.

3)  For all $e^{-1} \in E^{-1}$, $u^1(b^1, e^{-1}) \geq u^1(a^1, e^{-1})$.

*The action $b^i$ is called a temptation.*

In other words, an action is vulnerable if it is possible to lower the probability of all of the signals in $\hat{Y}$ by at least $\underline{\rho}$ while increasing the probability of each other signal by at least the multiple $(1 + \tilde{\rho})$. Notice that for an action to be vulnerable to a temptation, it must place at least weight $\underline{\rho}$ on each signal in $\hat{Y}$. Notice also that the definition does not control the payoff to the vulnerable action when the short-run players participate – the temptation here is not to increase short-run payoff, but rather to decrease the probability of the signals in $\hat{Y}$. In the EV example, the action *et* is vulnerable relative to $\{E\}$. The temptation $b^i$ is *tt*, which sends the probability of the signal $E$ to zero. (Since there is one other signal, condition 2 of the definition is immediate.)

Notice that if an action $a^1$ is vulnerable, it cannot be the case that if $\alpha^{-1} \notin \text{conhull } E^{-1}$ then $\rho(\cdot \mid a^1, \alpha^{-1}) = \rho(\cdot \mid \alpha^{-1})$ – the distribution of actions must be in some way dependent on the long-run player's action if the short-run players do not exit. This is related to the notion of an action being identified, as in Fudenberg, Levine and Maskin [1994]. Here we allow the possibility that there are strategies such as *et* and *te* from the EV example that are not identified, but do not allow complete lack of identification unless the short-run players play in $E^{-1}$ with probability one.

**Definition 6:** *A mixed action $\alpha^1$ for the long run player is enforceable if there does not exist another action $\tilde{\alpha}^1$ such that for all $a^{-1} \in E^{-1}$,*
$u^1(\tilde{\alpha}^1, a^{-1}) \geq u^1(\alpha^1, a^{-1})$ *and for all $a^{-1} \in A^{-1} - E^{-1}$,*

$u^1(\tilde{\alpha}^1, a^{-1}) > u^1(\alpha^1, a^{-1})$ *and* $\rho(\cdot \mid \tilde{\alpha}^1, a^{-1}) = \rho(\cdot \mid \alpha^1, a^{-1})$. *When* $\alpha^1$ *is not enforceable, we say that the action* $\tilde{\alpha}^1$ *defeats* $\alpha^1$.

If an action is not enforceable then there is necessarily lack of identification, since $\alpha^1$ and $\tilde{\alpha}^1$ induce exactly the same distribution over signals. The key point is that if the short-run players enter with positive probability, the rational type cannot play an action that is not enforceable: by switching to $\tilde{\alpha}^1$ he would strictly increase his current payoff, while maintaining the same distribution over signals, and so the same future utility. Note also that a mixed action that assigns positive probability to unenforceable actions is not enforceable: if $\alpha^1$ assigns probability $p$ to unenforceable action $a^1$, then $\alpha^1$ is defeated by the mixed action $\hat{\alpha}^1$ formed by replacing the probability on $a^1$ with the action $\tilde{\alpha}^1$ that defeats $a^1$.

**Definition 7:** *A participation game has an* exit minmax *if*

$$\max_{\alpha^{-1} \in E \cap range(B)} \max_{\alpha^1} u^1(\alpha^1, \alpha^{-1}) = \min_{\alpha^{-1} \in range(B)} \max_{\alpha^1} u^1(\alpha^1, \alpha^{-1})$$

In other words, any exit strategy forces the long-run player to the minmax payoff, where the relevant notion of minmax incorporates the restriction that the action profile chosen by the short-run players must lie in the range of $B$.[7] It is convenient in this case to normalize the minmax payoff to 0.

We are now in a position to define a class of games we call *bad reputation games*.

---

[7] When there is a single short-run player this restriction collapses to the constraint of not playing strictly dominated strategies, but when there are multiple short-run players it involves additional restrictions. It is clear that no equilibrium could give the long-run player a lower payoff than the minmax level defined in defintion 7. Conversely, in complete-information games, any long-run player payoff above this level can be supported by a perfect Bayesian equilibrium if actions are identified and the public observations have a "product structure" (Fudenberg and Levine [1994]). This is true in particular when actions are publicly observed as shown in Fudenberg, Kreps and Maskin [1990].

**Definition 8:** *A participation game is a* bad reputation game *if it has an exit minmax, there is an unfriendly set* $\widehat{A}^1$*, a friendly set* $F^1$ *that is orthogonal to* $\widehat{A}^1$*, and a set of signals* $\widehat{Y}$ *that are unambiguous for* $\widehat{A}^1$*, and such that every enforceable* $f^1 \in F^1$ *is vulnerable to temptation relative to* $\widehat{Y}$*.*

In particular, the EV game is a bad reputation game. We take the friendly set to be $\{et\}$, the unfriendly set to be $\{ee\}$ and the unfriendly signals to be $\{E\}$. We have already observed that $\{et\}$ is a friendly set and $\{ee\}$ unfriendly. The two are obviously orthogonal, and $\{E\}$ is unambiguous for $\{ee\}$.

In a bad reputation game, the relevant temptations are those relative to $\widehat{Y}$. For the remainder of the paper when we examine a bad reputation game and refer to a temptation, we will always mean relative to the set $\widehat{Y}$.

For any bad reputation game, it is useful to define several constants describing the game. Let $\widehat{a}^1$ be the probability in the definition of an unfriendly set; let $\underline{\alpha}$ be the probability in the definition of a friendly set. Since the friendly set is finite, we may define $\underline{\rho} > 0$ to be the least value for which a friendly enforceable action is vulnerable. . Define

$$r = \min_{\widehat{a}^1 \in \widehat{A}^1, a^1 \notin \widehat{A}^1, \alpha^{-1} \notin \text{conhull}(E^{-1}), \widehat{y} \in \widehat{Y}} \frac{\rho(\widehat{y} \mid \widehat{a}^1, \alpha^{-1})}{\rho(\widehat{y} \mid a^1, \alpha^{-1})}.$$

Since $\widehat{Y}$ is unambiguous for the unfriendly set, $r > 1$. Also define

$$\eta = -\log(\underline{\alpha}\underline{\rho}) / \log r$$

which is positive, and

$$k_0 = -\frac{\log(\widehat{\alpha})}{\log(\widehat{\alpha} + (1 - \widehat{\alpha})r)}.$$

## 3. The Theorem

We now prove our main result: In a bad reputation game with a sufficiently patient long-run player and likely enough unfriendly types, in

any Nash equilibrium, the long-run player gets approximately the exit payoff. The proof uses two Lemmas, both proven in the Appendix.

We begin by describing what it means for unfriendly types to be likely "enough." Let $\Theta(F^1)$ be the commitment types corresponding to actions in the support of $F^1$.

**Definition 9:** *A bad reputation game has* commitment size $\varepsilon, \omega, \phi$ *if*

$$\mu(0)[\Theta(F^1)] \leq \varepsilon \left( \frac{\varepsilon}{\omega} \frac{\mu(0)[\widehat{\Theta}]}{\mu(0)[\Theta(F^1)]} \right)^{\phi}$$

This notion of commitment size places a bound on the prior probability of friendly commitment types that depends on the prior probability of the unfriendly types in $\widehat{\Theta}$. Since $\phi$ is positive, the larger the prior probability of $\widehat{\Theta}$, the larger the probability of the friendly commitment types is allowed to be. The hypothesis that the priors have commitment size $\varepsilon, \widehat{\alpha}, \eta$ for sufficiently small $\varepsilon$ is a key assumption driving our main results.

Note that the assumption of a given commitment size does not place any restrictions on the relative probabilities of commitment types. In particular, let $\tilde{\mu}$ be a fixed prior distribution over the commitment types, and consider priors of the form $\lambda\tilde{\mu}$, where the remaining probability is assigned to the rational type. Then the right-hand side of the inequality defining commitment size depends only on $\tilde{\mu}$, and not on $\lambda$, while the left-hand side has the form $\lambda\tilde{\mu}$. Hence for sufficiently small $\lambda$ the assumption of commitment size $\varepsilon, \widehat{\alpha}, \eta$ is satisfied. Note that in EV the set of actions $\widehat{A}^1 = \{ee\}$ has commitment size $0, \widehat{\alpha}, \eta$ for all priors $\mu(0)$ since the only types are the rational type and the commitment type who plays $ee$.

We now have a series of Lemmas proven in the Appendix.

**Lemma 1:** *If $h_t$ is a positive probability history in which $\widehat{y} \in \widehat{Y}$ occurs in period $t$ and $\mu(h_{t-1})[\Theta(F^1)] \leq \underline{\alpha}/2$ then $\alpha_0^1(h_t) \geq (\underline{\alpha}/2) f^1$.*

In other words, when the prior on committed types is sufficiently low, entry can occur only if the strategic type is playing a friendly strategy with appreciable probability.

**Lemma 2:** *In a bad reputation game, if $h_t$ is a positive probability history with respect to a Nash equilibrium, and the signals in $h_t$ all lie in $Y^e \cup \widehat{Y}$, then*

*a) At most $k^* = k_0 - \log\left(\mu(0)[\widehat{\Theta}]\right)$ of the signals are in $\widehat{Y}$.*

*b) If the commitment size is $\underline{\alpha}/2, \widehat{\alpha}, \eta$ then $\mu(h_t)[\Theta(F^1)] \leq \underline{\alpha}/2$ whenever $\alpha^{-1}(h_t) \notin \text{conhull } E^{-1}$.*

*Remark*: The intuition for part *a* is simple, and closely related to the argument about the deterministic evolution of beliefs in FL: The short-run players exit if they think it is likely that entry will lead to the observation of a bad signal. Hence each observation of a bad signal is a "surprise" that increases the posterior probability of the bad type by (at least) a fixed ratio greater than 1, so along a history that consists of only bad signals and exit signals, the posterior probability of the bad type eventually gets high enough that all subsequent short-run players exit. This argument holds no matter what other types have positive probability, and it is the only part of this lemma that would be needed when there are only two types, one rational and one bad, as in EV.

However, as we will show by example below, we cannot expect the "bad reputation" result to hold when there is sufficiently high probability of the Stackelberg type. Part *b* of the lemma says that if the initial probability of the friendly types is sufficiently low compared to the prior probability of the bad types, then along any history consisting of exit outcomes and bad outcomes, the probability of the Stackelberg type remains low up to the point where the short-run players decide to exit. The intuition for this result is that because $r > 1$, each observation of a bad

signal not only increases the probability of the bad type, it increases the relative probability of this type compared to any friendly commitment type. (If there were a type with a history-dependent strategy, this part of the lemma would need to be modified.) Notice that part b will be satisfied for any given ratio of type probabilities, provided that the probability of all types is sufficiently low.

Define $U^1 = \max u^1 - \min(0, u^1)$,

$$\bar{u}^1(y, \tilde{\rho}) = \begin{cases} 0 & y \in Y^e \\ \left(1 + \dfrac{1}{\tilde{\rho}}\right)U^1 & y \in \hat{Y} \end{cases}$$

$$\bar{\delta}^1(y, \tilde{\rho}) = \begin{cases} \delta & y \in Y^e \\ \dfrac{\delta}{\tilde{\rho}} & y \in \hat{Y} \end{cases}$$

and $Y(h_t) = \{y \in Y^e \cup \hat{Y} \mid \rho(y \mid \bar{\alpha}^1(h_t), \alpha^{-1}(h_t)) > 0\}$.

**Lemma 3:** *In a participation game if $\alpha^{-1}(h_t) \in \mathrm{conhull}(E^{-1})$, or $\alpha^{-1}(h_t) \notin \mathrm{conhull}\, E^{-1}$ and $\alpha_0^1(h_t) \geq \gamma f^1$ for some $\gamma > 0$ and vulnerable friendly action $f^1$ of temptation $f$ size $\underline{\rho}, \tilde{\rho}$ then*

$$v^1(h_t) \leq \max_{y \in Y(h_t)}(1 - \delta)\bar{u}^1(y, \underline{\rho}) + \bar{\delta}(y, \tilde{\rho})v^1(h_t, y).$$

*Remark*: This lemma says that if the rational type is playing a friendly strategy, his payoff is bounded by a one-period gain and the continuation payoff conditional on a bad signal. This follows from the assumption that for every entry-inducing strategy it is possible to lower the probability of all of the signals in $\hat{Y}$ by at least $\underline{\rho}$ while increasing the probability of each other signal by at least the multiple $(1 + \tilde{\rho})$. The fact that the rational type chooses not to reduce the probability of the bad signal means that the continuation payoff after the bad signal cannot be much worse than the overall continuation payoff.

**Theorem 1:** *In a bad reputation game of commitment size $\underline{\alpha}/2, \hat{\alpha}, \eta$ let $\bar{v}^1$ be the supremum of all Nash equilibrium values for the rational type.*

$$\bar{v}^1 \leq (1 - \delta)k^* \left(\frac{1}{\tilde{\rho}}\right)^{k^*} \left(1 + \frac{1}{\tilde{\rho}}\right)U^1,$$

where $k^* = k_0 - \log\left(\mu(0)[\widehat{\Theta}]\right)$. In particular, $\lim_{\delta \to 1} \bar{v}^1 \leq 0$.

*Proof:* Given an equilibrium, we begin by constructing a positive probability sequence of histories beginning with 0. Given $h_t$ already constructed, we define $h_{t+1} = (h_t, y_{t+1})$ where

$$y_{t+1} \in \arg\max_{y \in Y(h_t)}(1 - \delta)\bar{u}^1(y, \underline{\rho}) + \bar{\delta}(y, \tilde{\rho})v^1(h_t, y).$$

We know that $Y(h_t)$ is not empty because either $\alpha^{-1}(h_t) \in$ conhull $E^{-1}$, or $\alpha^{-1}(h_t) \notin$ conhull $E^{-1}$. This latter case implies that $\bar{\alpha}^1(h_t) \geq \underline{\alpha}f^1$ for some friendly $f^1$, and since only enforceable actions can be played in equilibrium, this $f^1$ must be vulnerable to temptation, so $\rho(\widehat{Y} \mid \bar{\alpha}^1(h_t), \alpha^{-1}(h_t)) \geq \underline{\alpha}\rho(\widehat{Y} \mid f^1, \alpha^{-1}(h_t)) > 0$.

Now apply Lemma 2 to conclude that for each $h_t$ at most $k^*$ of the signals are in $\widehat{Y}$ and $\mu(h_t)[\Theta(F^1)] \leq \underline{\alpha}/2$ whenever $\alpha^{-1}(h_t) \notin$ conhull $E^{-1}$. Consider an $h_t$ such that $\alpha^{-1}(h_t) \notin$ conhull $E^{-1}$. From the definition of a friendly action, we know that $\bar{\alpha}^1(h_t) \geq \underline{\alpha}f^1$ for some friendly $f^1$, so $\mu(h_t)[\Theta(F^1)] \leq \underline{\alpha}/2$ and Lemma 1 implies that $\alpha_0^1(h_t) \geq (\underline{\alpha}/2)f^1$. Now apply Lemma 3 to conclude that for each $h_t$

$$v^1(h_t) \leq (1 - \delta)\bar{u}^1(y_{t+1}, \underline{\rho}) + \bar{\delta}(y_{t+1}, \tilde{\rho})v^1(h_{t+1}).$$

Since $v^1(h_t) \leq U^1$, it follows that

$$v^1(0) \leq (1 - \delta)\sum_{t=1}^{\infty} \Pi_{\tau=2}^t \bar{\delta}(y_\tau, \tilde{\rho})\bar{u}^1(y_t, \underline{\rho}).$$

Since $\bar{u}^1(y^e, \tilde{\rho}) = 0$, and $y_t \in \widehat{Y}$ at most $k^*$ times, this gives the desired bound. Notice that the fact that $\bar{u}^1(y^e, \tilde{\rho}) = 0$ follows from the assumption of exit minmax: it is here that we make use of the fact that exit gives the long-run player no more than the minmax.

☑

## 4. Examples

We now consider a number of examples to illustrate the scope of Theorem 1, and also the extent to which the assumptions are necessary as well as sufficient. To begin, Example 4.1 illustrates what happens when the prior puts too much weight on some committed types for the hypothesis of commitment size $\underline{\alpha}/2$ to be satisfied. Example 4.2 shows that the EV conclusion is not robust to the addition of an observed action that makes the short-run players want to enter. Example 4.3 examines participation games that are not bad reputation games, and example, 4.4 illustrates the role of the exit-minmax assumption. In all of the examples but 4.1, we assume that the hypothesis of commitment size $\underline{\alpha}/2$ is satisfied, and investigate whether the game is a bad reputation game. The following section considers a class of bad-reputation principal-agent games. .

### *Example 4.1: EV With Stackelberg Type*

We have verified Assumptions 1 and 2 in the EV example, so Theorem 1 follows. Moreover, we have relaxed the original assumptions of EV in a number of ways. One important extension is that we allow for positive probabilities of all commitment types. In particular, we allow a positive probability of a "Stackelberg type" committed to the honest strategy $et$, which is the optimal commitment. However, a hypothesis of the theorem is that the prior satisfy the commitment size assumption.

Here we illustrate that assumption in the context of the EV example. Suppose in particular that there are 3 types, rational, bad, and Stackelberg. The set of possible priors can be represented by the simplex in figure 2
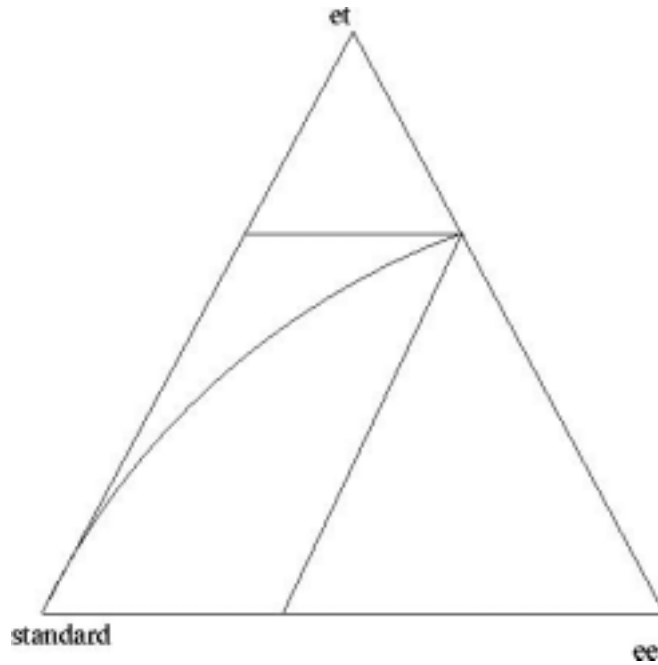


Figure 2

When the prior falls into the region in the lower right, the probability of the bad type is too high, and the short run players refuse to enter regardless of the behavior of the standard type. Bad reputation arises because the long-run player tries to prevent the posterior from moving into this region. In EV the prior assigned probability zero to the Stackelberg type. Thus the prior and all posteriors on the equilibrium path belong to the lower boundary of the simplex. When there is a sufficiently high probability of the Stackelberg type, the short-run players will enter regardless of the behavior of the standard type; this is the region at the top of the simplex. Note that the boundaries of these regions intersect on the right edge of the simplex: this point represents the mixture between *ee* and *et* which makes the short-run player indifferent between entry and exit.

When the prior falls in the bad region, there will be no entry and the long-run player obtains the minmax payoff of zero. On the other hand, when the prior falls in the good region, we there is a Nash (and indeed sequential) equilibrium in which the long run player receives the e best commitment payoff, which is "$u$" in the notation of EV.

Consider the game in which the posterior probability of the bad type is zero. In this game there exists a sequential equilibrium in which the long-run player gets $u$. Suppose that we assume that this is the continuation payoff in the original game in any subform in which the long-run player played $t$ at least once in the past. A sequential equilibrium of this modified game is clearly a sequential equilibrium of the original game, and by standard arguments, this modified game has a sequential equilibrium. How much does the rational long-run player get in this sequential equilibrium? One option is to play $tt$ in the first period. Since the short-run player is entering regardless, this means that beginning in period 2 the rational type gets $u$. In the first period he gets $(u - w)/2$. Hence in equilibrium he gets at least $(1 - \delta)(u - w)/2 + \delta u$, which converges to $u$ as $\delta \to 1$.

Our theorem is about the set of equilibrium payoffs for priors outside of these two regions. The theorem states that there is a curve, whose shape is represented in the figure, such that when the prior falls below this curve, the set of equilibrium payoffs for the long-run player is bounded above by a value that approaches the minmax value as the discount factor converges to 1. The diagram shows that the left boundary of the simplex is an asymptote for this curve as it approaches the complete information prior (i.e. $\mu(0)(\theta_0) = 1$) in the lower left corner. This illustrates the important aspect of the commitment size restriction: it is satisfied for r "almost) all" sufficiently small perturbations of the complete information limit.

### *Example 4.2: Adding an Observed Action to EV*

We now modify the EV game by adding new observable action "*g*" for the long-run player called "give away money." This action induces the short-run players to participate ($B(g) \cap \text{conhull}(E^{-1}) = \varnothing$) and it is observable, so Assumption 1 is not satisfied. Moreover, even without a Stackelberg type the EV conclusion fails in this game: there is an equilibrium where the rational type plays $g$ in the first period. This reveals that he is the rational type, and there is entry in all subsequent periods, while playing anything else reveals him to be the bad type so that all subsequent short run players exit. Thus the assumption that every friendly action is vulnerable to temptation is seen to be both important and economically restrictive.

### *4.3. Orthogonality Issues*

Suppose friendly actions send the bad signal by putting positive weight on unfriendly actions. An important class of games in which this is the case are those in which, conditional on entry, the long-run players' actions are observed. In this case the bad signals correspond to unfriendly actions, and bad signals can only have positive probability when the unfriendly action is played with positive probability. Moreover, in some games, the only friendly strategies involve randomizing in this way.

**Proposition 1:** *If there is a friendly action that only send bad signals because of mixing onto unfriendly actions, the game is not a bad reputation game*

Proof: The assumption that the friendly and unfriendly sets are orthogonal are violated.

☑

To see that this makes a difference, consider the following two-person game:

|   | L | M | R |
|---|---|---|---|
| U | 0,4 | 1,3 | 0,0 |
| D | 0,0 | 1,3 | 0,4 |

Figure 3

where $L$ and $R$ correspond to exit and $M$ to entry.[8] In this case entry can be induced only by mixing with probability of $U$ between ¼ and ¾. Because friendly actions must involve mixing, they will send a bad signal, which can be taken to be either $U$ or $D$. Suppose we take $D$ to be the bad signal, and suppose that the only committed type with positive probability is $D$. The problem that occurs is that while repeated observations of $D$ increase the probability of the bad type, when that probability hits the critical level, the rational player no longer needs to play a friendly action: in effect the bad type is doing his mixing for him. Specifically, suppose that initially, the probability of the bad type is less than ¼, and that for any current probability $\mu(h_t)[D]$ less than ¼ the rational type mixes so that the overall probability of $D$ is exactly ¾. The short-run player always enters. If $U$ is observed, the type is revealed rational. If $D$ is observed, the probability of the bad type increases by a factor of $4/3$. So when it first exceeds ¼ it is at most equal to 1/3. At this point, the rational type may reveal himself by playing $U$ with probability 1, while preserving the incentive of the short-run player to enter. In this equilibrium, the long-run player gets 1.

We say that an action $f^*$ is *sufficient for entry* if, for some $\underline{\alpha} < 1$, $\alpha^1 \geq \underline{\alpha} f^*$ implies that there is $\alpha^2 \in B(\alpha^1)$ with positive probability of entry. In the example above the friendly action is sufficient for entry, and sends the bad signal only because of mixing onto the

---

[8] In this example, the short-run player has several exit actions, and his payoff depends on the long-run player's action. This is a necessary feature of two-player games where the only friendly strategies are mixed, but it is not necessary in three-player games – think of a game where player 3 has veto power, 3 only plays In if 2 plays M, and 2's payoffs to M are as in the payoff matrix of this example.

unfriendly action. That is, the sufficient action mixes between a pure action that does not send the bad signal, and an unfriendly action. If there is a friendly action that does not send the bad signal at all, then we have a quite general conclusion that the game is not a bad reputation game since such an action cannot admit a temptation. More strongly, if an action sufficient for entry does not send the bad signal at all, then a patient rational player can do almost as well as in the absence of bad reputation effects.

**Proposition 2:** *If there is an action $f^*$ that is sufficient for entry and does not send any bad signal, the only committed types are unfriendly types, and the probability of committed types is sufficiently low, then as $\delta \to 1$ there are sequential equilibrium payoffs for the rational type that approach the highest sequential equilibrium payoff without committed types.*

*Proof:* Suppose that the prior probability of committed types is sufficiently low that the short-run players will enter when the rational type plays $f^*$. Then it is a sequential equilibrium for the rational type to play $f^*$ in the initial period with entry by the short-run players. Subsequently, if a bad signal was observed, the short-run players stay out. If a bad signal was not observed, the probability of committed types is zero, and the continuation equilibrium is the best possible without committed types. On the equilibrium path, the rational type payoff clearly approaches that of the highest payoff without committed types, since he gets that amount beginning in period 2, and payoffs in period 1 are bounded below.

☑

### 4.4: Exit Minmax

In participation games, reputation plays a role because the short run players will guard against unfriendly types by exiting. This is "bad" for the long-run player only if exit is worse than the payoff he otherwise

would receive, and the exit minmax assumption ensures that this is the case.

In participation games without exit minmax, there are outcomes that are even worse for the long-run player than obtaining a bad reputation. In this case it is possible that there exist equilibria in which the long-run player is deterred from his temptation to avoid exit by the even stronger threat of a minmaxing punishment. For example consider the game in Figure 4, where the first matrix represents the payoffs, and the second represents the distribution of signals conditional on entry.

| | $In$ | $Out^1$ | $Out^2$ |
|---|---|---|---|
| $F$ | 1,1 | 0,0 | -2,0 |
| $U$ | 1,0 | 0,1 | -2,0 |
| $T$ | 1,0 | 0,0 | -2,1 |

| | $g$ | $b$ | $r$ |
|---|---|---|---|
| $F$ | ½ | ½ | 0 |
| $U$ | 0 | 1 | 0 |
| $T$ | ½ | 0 | ½ |

Figure 4

This game is a participation game with exit actions $Out_1$ and $Out_2$, unfriendly action $U$ and friendly action $F$ vulnerable to temptation $T$. There are only two types, the rational type and a bad type that plays $U$. Exit minmax fails because the maximum exit payoff exceeds the minmax payoff, and we claim that there are good equilibria in this game because the threat of exiting with $Out_2$ is worse than the fear of obtaining a reputation for playing $U$ which would only lead to exit with $Out_1$.

To see this, consider the following strategy profile. The rational type plays $F$ at every history unless the signal $r$ has appeared at least once; in that case the rational type plays $T$. The short run player plays $Out_2$ if a signal of $r$ has ever appeared. Otherwise, the short run player plays $Out_1$ if the posterior probability of the bad type exceeds ½ and $In$ if this probability is less than ½. Obervations of $r$ are interpreted as signals that the long-run player is rational.

Since $(T, Out_2)$ is a Nash equilibrium of the stage game, the continuation play after a signal of $r$ is a sequential equilibrium. When $r$

has not appeared, the long run player optimally plays $F$. Playing $U$ gives no short-run gain and hastens the onset of $Out_1$, and playing $T$ shifts reallocates probability from the bad signal $b$ to the signal $r$ which is even worse.[9] The short-run players are playing short-run best responses. In this equilibrium, the long run player does not give in to the temptation to play $T$. As a result, with positive probability, the short-run players never become sufficiently pessimistic to begin exiting, and so the long run player achieves his best payoff.

In the above example there were two exit actions. The next proposition states that when there is only one exit action and the long-run player's exit payoff is independent of his own action, the worst Nash equilibrium payoff for the long run player is (not much worse than) his exit payoff. Note that this condition is satisfied in the principal-agent applications discussed in section 5. The proposition is a consequence of FL (1992).

**Proposition 3:** *Consider a participation game with a single short-run player and a unique exit action. If*

*(i) there exists a pure action[10] $\hat{a}^1$, such that $B(\alpha^1) = \{exit\}$,*

*(ii) the prior distribution assigns positive probability to a type that is commited to $\hat{a}^1$,*

and

(iii) the long-run player's action is identified conditional on entry,

*then there is a lower bound on the payoffs to the standard type which converges to $\bar{u}$ as the discount factor approaches 1.*

---

[9] Playing $T$ gives probability ½ of shifting to the absorbing state where payoffs are –2. Playing the equilibrium action of $F$ has probability at most ½ of switching to the state where payoffs are 0.

[10] The assumption of pure action is not needed here, but we state the result this way for consistency with the rest of the paper.

*Proof:* FL (1992) established[11] that for any game there exists a bound lower bound $b(\delta)$ on the set of Nash equilibrium payoffs for the standard type, and that as $\delta \to 1$, $b(\delta)$ converges to a limit that is at least

$$\max_{a^1 \in C} \min_{\alpha^2 \in \tilde{B}(a^1)} u^1(\alpha^1, \alpha^2)$$

where $\tilde{B}(a^1)$ is the set of self-confirmed best-responses to for the short-run player to $a^1$, and $C$ is the set of actions corresponding to the support of the prior distribution over commitment types. Because the long-run player's action is identified conditional on entry and $B(\hat{a}^1) = \{exit\}$, we have $\tilde{B}(\hat{a}^1) = \{exit\}$, and because the type that plays $\hat{a}^1$ has positive prior probability, the FL (1992) bound is at least $u^1(\hat{a}^1, exit) = \bar{u}$ .

☑

For games satisfying the conditions of the proposition, the exit minmax condition is not necessary for bad reputation. The worst equilibrium continuation value that the short-run players could inflict is arbitrarily close to the exit payoff and hence a patient long run player could not be deterred from his temptation to avoid a bad reputation.

## 5. Poor Reputation Games and Strong Temptations

Recall that an action is vulnerable to a temptation if when the short-run players participate, the temptation lowers the probability of all bad signals, and increases the probability of all others. In this case the bad reputation result requires the exit minmax condition, as demonstrated by the example in Section 4.4. Notice, however, that in the example the relative probability of $g$ and $r$ is changed by the temptation. If the temptation satisfies the stronger property that the relative probability of the other signals remains constant, then we can weaken the assumption of

---

[11] The statement of the FL theorem requires that commitment types including mixing types have full support, in which case the set $C$ is the space of all (mixed) actions, but the proof given there also shows that the version of the lower bound given here is correct.

exit minmax. In this section we prove this result, and give an application to games with two actions.

First we give a formal definition of a strong temptation:

**Definition 10:** *An action* $a^1$ *is* vulnerable to a strong temptation relative to *a set of signals* $\widehat{Y}$ *if there exists a number* $\underline{\rho} > 0$ *and an action* $b^i$ *such that*

*1) If* $a^{-1} \notin E^{-1}$, $\widehat{y} \in \widehat{Y}$ *then* $\rho(\widehat{y} \,|\, b^1, a^{-1}) \leq \rho(\widehat{y} \,|\, a^1, a^{-1}) - \underline{\rho}$

*2) If* $a^{-1} \notin E^{-1}$ *and* $y, y' \notin \widehat{Y} \cup Y^e$ *then* $\dfrac{\rho(y \,|\, b^1, a^{-1})}{\rho(y' \,|\, b^1, a^{-1})} = \dfrac{\rho(y \,|\, a^1, a^{-1})}{\rho(y' \,|\, a^1, a^{-1})}$.

*3) For all* $e^{-1} \in E^{-1}$, $u^1(b^1, e^{-1}) \geq u^1(a^1, e^{-1})$.

The action $b^i$ is called a strong temptation.

The first and third parts of this definition are the same as in the definition of a temptation; the additional strength comes from part (2), which requires that the temptation not merely increase the probability of all of the good signals, but leave their relative probabilities unchanged. Note that strong temptation is equivalent to temptation in games in which the set $Y \setminus (\widehat{Y} \cup Y^e)$ has a single element, for example games in which there are only two entry signals; in particular applies when the modified versions of the game of Section 4.4 is modified so that the only signals when entry occurs are *g* and *r*.

This condition lets us sharpen lemma 3 by replacing the variable $\overline{\delta}^1(y, \tilde{\rho})$ with the constant $\delta$:

**Lemma 4:** *In a participation game, if* $\alpha^{-1}(h_t) \in \text{conhull}\, E^{-1}$ *or* $\alpha^{-1}(h_t) \notin \text{conhull}\, E^{-1}$ *and* $\alpha_0^1(h_t) \geq \gamma f^1$ *for some* $\gamma > 0$ *and friendly action* $f^1$ *that is vulnerable to a strong temptation size* $\underline{\rho}$, *then*

$$v^1(h_t) \leq \max_{y \in Y(h_t)} (1 - \delta) \overline{u}^1(y, \underline{\rho}) + \delta v^1(h_t, y).$$

The proof, in the Appendix, follows that of lemma 3, but takes advantage of the fact that the long-run player's continuation expected value, conditional on a friendly action, a non-exit profile, and a signal not in $\widehat{Y} \cup Y^e$, is the same for the equilibrium action and the strong temptation $b^1$.

Define

$$\hat{u}^1 = \max_{a^1, \alpha^{-1} \in \text{conhull}(E^{-1}) \cap \text{image}(B)} u^1(a^1, \alpha^{-1})$$

This is a bound on the long-run player's payoff when the short-run players play exit actions that are a best response to some (possibly incorrect) conjectures.

**Defintion 11:** *A participation game is a* poor reputation game *if there is an unfriendly set $\widehat{A}^1$, a friendly set $F^1$ that is orthogonal to $\widehat{A}^1$, and a set of signals $\widehat{Y}$ that are unambiguous for $\widehat{A}^1$, and such that every enforceable $f^1 \in F^1$ is vulnerable to strong temptation relative to $\widehat{Y}$.*

**Theorem 3:** *In a poor reputation game of commitment size $\underline{\alpha}/2, \widehat{\alpha}, \eta$*

$$\lim_{\delta \to 1} \overline{v}^1 \leq \hat{u}_1.$$

In other words, poor reputation games have much the same consequences as bad reputation games. Notice that it is possible for a game to be both a bad reputation game and a poor reputation game, and, since strong and ordinary temptation are equivalent when $Y \setminus (\widehat{Y} \cup Y^e)$ the two are necessarily equivalent in this case. The original EV game is such an example. Notice also that example 4.4 in which we construct a non-bad equilibrium has three signals rather than two. With two signals, the game would still fail the exit minmax condition and fail to be a bad reputation game, but it would never-the-less be a poor reputation game, and would not admit a good equilibrium. Finally, observe that there is an element of continuity: the proof of both Lemma 3 and 4 can be generalized, so that the extent to which the best equilibrium (in the limit as $\delta \to 1$) can exceed the most favorable outcome with exit is bounded by a term which

is the product of the change in relative probabilities induced by a temptation and the excess of the best result given exit over the minmax. When one of these two equals zero we get the case of either bad or poor reputation. Otherwise, the best equilibrium can exceed the best exit payoff for the long-run player, but only by a limited amount.

We turn now to the special case of two-player participation games where there is only one signal in $\widehat{Y}$ and short-run player payoffs depend only on the signal. We focus on the case where bad reputation games have poor reputation, that is one signal in $Y \setminus (\widehat{Y} \cup Y^e)$. We show that these games are not poor reputation games (and by implication not bad reputation games either)..

***Proposition 4:*** In a two-player participation game suppose there are only two "entry signals" (that is two elements of $Y - Y^e$), that the short-run player has only two actions, and that the short-run player's realized payoff is determined by the signal. Then the game is not a poor reputation game. *Proof:* Notice that since the short-run player has only two actions, they correspond to "entry" and "exit" respectively. Consequently, the short-run player payoff conditional on entry depends only on the distribution over signals induced by the long-run player action. If we normalize the short-run player's payoff function so that his exit payoff is 0, and suppose that both the friendly and unfriendly sets are non-empty, then one signal yields a negative payoff and the other signal's payoff is positive; call these the "bad" and "good" signals respectively. Then any unfriendly set $\widehat{A}^1$ consists of actions with a sufficiently high probability of sending the bad signal, and the bad signal (as a singleton set) is the only set $\widehat{Y}$ that can be unambiguous for $\widehat{A}^1$. (If there were no unfriendly set, then the game is not a poor reputation game, so we can assume there is at least one unfriendly set.) Let $f^1$ be the friendly action in the (finite) friendly set that maximizes the short-run player's payoff. The payoff to this action, conditional on it not generating the bad signal with the negative payoff, is positive, and since any temptation relative to $\widehat{Y}$ must reduce the probability of the bad signal, a temptation must give the short-

run player a higher payoff than this "friendliest" friendly action. For this to be true, there must be a pure strategy $\hat{b}^1$ that gives the short-run player at least this same utility. Clearly $\hat{b}^1$ induces entry, and since it is a pure strategy, it must be in the friendly set. This contradicts the fact that $f^1$ was assumed to maximize short-run player utility in the friendly set.

<div align="right">☑</div>

We believe that the assumptions of this proposition imply that there is an equilibrium where the rational type's payoff is bounded below by a positive number as $\delta \rightarrow 1$ but we have not been able to show this.

## 6. Principal-Agent Entry Games

In this section we consider a class of applications which have the nature of an agency relationship. The long- run player (the agent) takes an action that affects the payoffs of both a principal (that period's short run player) and herself. When the principal's and the agent's preferences differ over the action set, and the action is not perfectly observed, we have a classical problem of incentives. A repeated interaction can often substitute for explicit contracts in alleviating this incentive problem. The long run agent's objective of establishing a good reputation can provide an incentive for efficient behavior in the short-run. In this section we classify agency environments in which the repeated interaction has the opposite effect. Bad reputation can intensify rather than mitigate the agent's incentive problem.

There is a single short-run player (the principal) whose only choice is whether to enter or to exit. If the principal enters, then the long-run player (the agent) chooses a payoff-relevant action, otherwise both players receive a reservation value which is normalized to zero. Formally $A^2 = \{exit, enter\}$ and $u^2(a^1, exit) = 0$ for each $a^1 \in A^1$. For simplicity we write $u^2(a^1, enter) = u^2(a^1)$. We assume there is an action $a^1 \in A^1$ for which $u^1(a^1) \geq 0$, so that the exit minmax assumption is satisfied. (Note that this assumption will hold whenever the principal has the option to

refuse to participate. Note also that from Theorem 2 this assumption is not necessary for games with two signals.)

For these games we can immediately identify the relevant friendly set. Define

$$F^1 = \{a^1 \in A^1 : u^2(a^1) \geq 0\}$$

which is the set of pure friendly actions. We know that $F^1 \subset \hat{F}^1$ for any friendly set $\hat{F}^1$. In fact, within the class of principal-agent games, any bad reputation game is a bad reputation game with friendly set $F^1$. To see this note that if $\alpha^1(F^1) = 0$ then $u^2(\alpha^1) < 0$, i.e. exit is the unique best-reply to $\alpha^1$. Thus $F^1$ is itself a friendly set.[12] Furthermore $\text{supp}(F^1) \subset \text{supp}(\hat{F}^1)$ so that orthogonality is preserved, and if every $f^1 \in \hat{F}^1$ is vulnerable then every $f^1 \in F^1$ is vulnerable. Thus we can restrict attention to $F^1$.

To show that these are bad reputation games, it suffices to find an unfriendly set orthogonal to $F^1$ with unambiguous signals, such that every enforceable point in $F^1$ is vulnerable to a temptation.

*Remark:* It is also of interest to consider games in which the agent has the opportunity to take a costly action prior to the entry decision of the short-run players. Consider for example, a game in which the long-run player is an expert advisor, and the decision of the short-run player is whether or not to pay the long-run player for advice. One example of this is the EV example of car repairs, where the long-run player is able to determine the type of repair the car needs. Other examples include stockbrokers advising clients on portfolio choices, doctors advising patients on treatments, and the IMF advising countries on economic policies. Costs incurred on exit are consistent with a bad reputation game provided that conditional on exit the temptations are less costly than the friendly actions. For example, the long-run player might be a stockbroker, and the general non-client specific information might be something about general economic conditions,

---

[12] When there is more than one principal, this conclusion does not follow, and mixed friendly sets will generally have to be considered. See the discussion in section 4.3 and footnote 7.

acquired in advance in the form of economic reports that will be presented to the client. The friendly actions in this case are to report truthfully; the bad action might be to always claim that times are good. In this case the temptation is to announce that times are bad when they are actually good, to avoid being mistaken for the type that always announces good times. If it is costly to put together a persuasive package of economic data indicating that times are bad when in fact they are good this would not be a bad reputation game. If it is more costly to put together an honest report, then it would be a bad reputation game.

### *6.1 Games with Hidden Information*

In these games the principal has some private information that is relevant for a decision affecting both principal and agent. Each period, nature draws a state $\omega \in \Omega$; in independently from a probability distribution that we denote by $p$.[13] The agent privately observes the state and then selects a decision $d \in D$. Conditional on the realized state and the decision of the agent, a signal $z \in Z$, is drawn from the distribution $m(z \mid \omega, d)$ where we assume that $m(z \mid \omega, d) > 0$ for each $z, \omega$, and $d$. Future short run players observe both $z$ and the decision $d$. Each player $j$ has state-dependent utility function $\pi^j(\omega, d)$ and evaluates stage payoffs according to expected utility with respect to $p(\omega)$.

To apply Theorem 1, we find conditions under which this defines a bad reputation game. The set of actions $A^1$ for the long-run player is the set of maps $a^1 : \Omega \to D$. The stage-game utility function is

$$u^j(a^1) = \sum_{\omega \in \Omega} p(\omega)\pi^j(\omega, a^1(\omega)).$$

Finally $Y - Y^e = Z \times D$ and

$$\rho((z,d) \mid a^1, entry) = \sum_{\{\omega : a^1(\omega) = d\}} p(\omega)m(z \mid \omega, d)$$

---

[13] This is a slight abuse of notation, as $p$ also denotes the probability distribution over types in the incomplete-information games, but no ambiguity should result.

**Proposition 5:** *The hidden information game is a bad reputation game if there exists a decision $d$ such that $a^1 \in F^1$ implies $\varnothing \neq \{\omega : a^1(\omega) = d\} \neq \Omega$.*

*Proof:* Let $a(d)$ denote the constant action that chooses $d$ regardless of the signal $\omega$, and take $\widehat{A} = \{a(d)\}$. Because $m(z \mid \omega, d) > 0$, the set of signals $Y^d = Z \times \{d\}$ is unambiguous for $\widehat{A}$. If $a^1$ is friendly, then $\Omega_d = \{\omega : a^1(\omega) = d\} \neq \varnothing$. For each $b \neq d$ let $b^1$ be the action defined by $b^1(\omega) \neq d, \omega \in \Omega_d$ and $b^1(\omega) = a^1(\omega), \omega \notin \Omega_d$. Then $a(d)$ is vulnerable to any mixed strategy that puts positive weight on every $b^1$.

☑

Many examples can be found that meet the condition of the proposition. First of all, note that the EV example is a special case. In fact the theorem extends the example to allow for public signals $z$ about the short run players' realized payoffs (which are determined by $(\omega, d)$.

### 6. 2. Games with Hidden Actions

On the other hand, agency games with hidden actions, or moral hazard, tend not to be susceptible to bad reputation effects. The problem is that the second part of the definition of temptation typically fails because deviations will generally lower the probability of some good signals. However, a special case in which a hidden action game is a bad reputation game occurs when there is only one short-run player and only two signals.

The following proposition is an immediate application of the definition of a bad reputation game in this setting.

**Proposition 6:** *Suppose in a principal-agent entry game that $Y - Y^e = \{y^L, y^H\}$ and that $\widehat{a}^1$ strictly maximizes the probability of $y^L$ with $u^2(\widehat{a}^1) < 0$. If for every friendly enforceable $a^1$ there is a $b^1$ such that $\rho(y^L \mid b^1) < \rho(y^L \mid a^1)$ the game is a bad reputation game.*

We consider two applications of this idea. In the first, the agent chooses an action from a one-dimensional set ordered so that higher

actions are more likely to give rise to the high signal. Specifically, we let $A^1 = \{\underline{a}^1,...,\overline{a}^1\} \subset \mathfrak{R}$ and $\rho(y^H | a^1)$ are an increasing function of $a^1$. We assume that $u^2(a^1)$ is concave so that $F^1$ is an interval. Whether or not the game is a bad reputation game then depends on whether the principal prefers extreme or interior actions.

**Proposition 7:** *The hidden action game with two non-exit signals is a bad reputation game if and only if $\{\underline{a}^1, \overline{a}^1\} \subset A^1 - F^1$.*

*Proof:* Suppose $\{\underline{a}^1, \overline{a}^1\} \subset A^1 - F^1$. Then $\hat{a}^1 = \underline{a}^1$ and every friendly action $a^1 > \hat{a}^1$ is vulnerable to the (unfriendly) temptation $\overline{a}^1$. On the other hand if, $\overline{a}^1 \in F^1$, then the only candidate set of bad signals is $\hat{Y} = \{y^L\}$, meaning that $\overline{a}^1$ is not vulnerable to a temptation. In case $\underline{a}^1 \in F^1$, we simply reverse the role of the signals.

$$\boxtimes$$

In these two-signal games, as in the hidden information games, short-run player utility depends on aspects of the long-run player strategy that is unobserved by subsequent short-run players. Proposition 4 shows that this must be the case for a game with two entry signals to be a bad-reputation game.

### *6.3. Rules vs. Discretion*

We can build on the analysis of hidden information games to discuss the emergence of rules over discretion in agency relationships. To motivate the idea, consider college admissions. The university (the long-run agent) receives an application. The applicant is described by a set of characteristics $\omega \in \Omega = \Omega^o \times \Omega^n$. Some ($\Omega^o$) of these characteristics are publicly observable (for example race and SAT scores) and others ($\Omega^n$) are observed only by the university. This may include information that is truly private (like an interview) or information that require the expertise of the agent to interpret (for example, the strength of the applicant's high school.) A pure strategy for the university is a map from characteristics to the decision space $D=$ (admit, deny). The probability of

drawing characteristics $\omega$ is $p(\omega) > 0$. The university's preferences over applicants are summarized by the payoff function $\pi^1(\omega)$ if the student is admitted, and $R$ if the student is denied.

The short-run principal, player 2, is the state governor who chooses between allowing the university discretion in admissions, or imposing a rigid admission rule based on observable characteristics. There are many possible rules that the principal might use, but since she is a short-run player we can restrict attention to the rule that maximizes the principal's expected short-run payoff. This rule is a mapping $g : \Omega^o \rightarrow D$ that mandates admission if and only if $\omega^o \in g^{-1}(admit)$. The imposition of a rigid admission rule represents "exit." The public signal at date $t$ is $y_t = (d_t, a_t^2, \omega^o)$, where any signal with $a_t^2 = rule$ is an exit signal. The governor shares the same preferences as the university, receiving a utility of $\pi^1(\omega)$ for admits and $R$ for rejects.

Because the university can always implement $r$ on its own, exit minmax condition is satisfied. In order for discretion to improve upon $r$, for some set of verifiable characteristics, the admission decision should depend on the unverifiable characteristics. That is $a^1 \in F^1$ only if $a^1(\omega^o, \omega^n) = admit$ and $a^1(\omega^o, \hat{\omega}^n) = deny$ for some $\omega^n, \hat{\omega}^n$ and $\omega^o$. Then by essentially the same argument as in the hidden information case, the game is a bad reputation game with unfriendly set

$$\{a^1 : a^1(\omega^o, \cdot) = deny\}$$

For example, $\omega^o$ may be racial characteristics, and this unfriendly set represents the governor's fear that the university admissions are biased against members of the race in question

## 7. Mulilateral Entry Games

We now consider games with multiple principals. In these "mutilateral entry" games, in the short-run players choose only whether to participate or exit. If any short-run player chooses to exit, that player receives the reservation payoff of 0, but play between the agent and other principals is unaffected. That is, $A^j = \{exit, enter\}$ for each $j > 1$, and the

unique exit profile is $a_e^{-1} \equiv (exit,...,exit)$. The payoff of the short-run players who enter depends only on the action of the principal, and not on how many other short-run players chose to enter; to simplify notation we denote this "entry payoff" as $u^j(a^1)$. If all principals exit, the long-run player's payoff is 0; if $m$ of them choose to enter, the long-run player's payoff is $u^1(a^1,m)$. We assume that the agent cannot be forced to participate, so that there exists an action $a^1$ such that for all m, $u^1(a^1,m) \geq 0$.

We do not require that $u^1(a^1,m)$ is linear in $m$, so this class of games includes those in which the agent has the opportunity to take a costly action prior to the entry decision of the short-run players. Consider for example, a game in which the long-run player is an expert advisor, and the decision of the short-run player is whether or not to pay the long-run player for advice. One example of this is the EV example of car repairs, where the long-run player is able to determine the type of repair the car needs. Other examples include stockbrokers advising clients on portfolio choices, doctors advising patients on treatments, and the IMF advising countries on economic policies. In the EV example, the private information emerges as a consequence of the decision of the short-run player to consult the long-run player, so the advice is specific to the short-run player. In another cases, at least some part of the information is not specific to the short-run player. The advisor receives a report about the general desirability of various actions, and then meets with each of his $n$ short-run customers, possibly learning about their individual needs. Here the advisor receives the signal regardless of whether or not he is consulted by any particular short-run player, and he may incur costs ahead of time for doing so. That is, the long-run player's payoff may depend on his action even if the short-run players decline to participate.

Costs incurred on exit are consistent with a bad reputation game provided that conditional on exit the temptations are less costly than the friendly actions. For example, the long-run player might be a stockbroker, and the general non-client specific information might be something about

general economic conditions, acquired in advance in the form of economic reports that will be presented to the client. The friendly actions in this case are to report truthfully; the bad action might be to always claim that times are good. In this case the temptation is to announce that times are bad when they are actually good, to avoid being mistaken for the type that always announces good times. If it is costly to put together a persuasive package of economic data indicating that times are bad when in fact they are good this would not be a bad reputation game. If it is more costly to put together an honest report, then it would be a candidate for a bad reputation game.

We have the following obvious extension of Proposition 5.

**Proposition 8:** *Suppose in a multilateral entry game that* $Y - Y^e = \{y^L, y^H\}$ *and that* $\hat{a}^1$ *strictly maximizes the probability of* $y^L$ *with* $u^j(\hat{a}^1) < 0$. *If for every friendly enforceable* $a^1$ *there is a* $b^1$ *such that* $\rho(y^L \mid b^1) < \rho(y^L \mid a^1)$ *the game is a bad reputation game.*

For concreteness, suppose that the short-run players are students, the long-run player a teacher, and the signals are teaching evaluations. (This model could apply equally well to the decision to attend a particular college, graduate school, or take a particular job.)  Each period, each short-run player decides whether to enter - that is, take the class, or not. The long run player has a pair of binary choices: he can either teach well or teach poorly, and he can either administer teaching evaluations honestly or manipulate them. The public signals are whether the evaluations (averaged over respondents) are good, $y^H$ or poor, $y^L$. If the evaluations are administered honestly and the class is taught well, there is probability .9 of a good evaluation. If evaluations are administered honestly and the class is taught poorly, the probability of good evaluations is only .1. Manipulating the evaluations is certain to lead to a good evaluation.  All players get 0 if no students decide not to take the class. For a short-run player who enters, the short run player's payoffs are +1 for good teaching and -1 for bad. Let $m$ denote the number of students who take the class. The rational type of

long-run player pays a cost of $m$ to teach well; good evaluations are worth $2m$, while manipulating evaluations costs $3m$. Hence in the one-shot game with only the rational type, the unique sequential equilibrium is for the rational type to teach well and not manipulate the evaluations, for an expected payoff of .8.

However, when there is a small probability that the instructor is a bad type, and the instructor faces a sequence of short-run students, Proposition 7 applies. To see this, we see that teaching poorly and administering the evaluations honestly is the unfriendly action $\hat{a}^1$. The friendly set consists of the pure actions "teach well, administer honest evaluations" and "teach well, manipulate." Crucially, the action "teach well, manipulate" is unenforceable: teach poorly and manipulate yields a higher stage game payoff and the same distribution over signals. So the only enforceable action in the friendly set is "teach well, administer honestly." This admits the temptation "teach poorly, manipulate." Here the short-run player recognizes that if the long-run player chooses not to send the signal honestly, he loses his incentive to teach well, and so there is no reason to enter

## Appendix: Proofs

**Lemma 1:** *If $h_t$ is a positive probability history in which $\hat{y} \in \hat{Y}$ occurs in period $t$ and $\mu(h_{t-1})[\Theta(F^1)] \leq \underline{\alpha}/2$ then $\alpha_0^1(h_t) \geq (\underline{\alpha}/2)f^1$.*

*Proof:* Given $h_{t-1}$ the short-run players' profile has positive probability on a profile that does not exit. At such profiles $\bar{\alpha}^1(h_t) \geq \underline{\alpha}f^1$ for some friendly $f^1$. Since $\mu(h_{t-1})[\Theta(F^1)] \leq \underline{\alpha}/2$ friendly and unfriendly sets are orthogonal we see that $\alpha_0^1(h_t) \geq (\underline{\alpha}/2)f^1$.

☑

**Lemma 2:** *In a bad reputation game, if $h_t$ is a positive probability history with respect to a Nash equilibrium, and the signals in $h_t$ all lie in $Y^e \cup \hat{Y}$*

*a) At most* $k^* = k_0 - \log\big(\mu(0)[\widehat{\Theta}]\big)$ *of the signals are in* $\widehat{Y}$.

*b) If the commitment size is* $\underline{\alpha}/2, \widehat{\alpha}, \eta$ *then* $\mu(h_t)[\Theta(F^1)] \le \underline{\alpha}/2$
*whenever* $\alpha^{-1}(h_t) \notin$ conhull $E^{-1}$.

*Proof:* First observe that if $\mu(h_{t-1})[\widehat{\Theta}] \ge \widehat{\alpha}$, then the short-run players must exit in period $t$, so $\mu(h_t) = \mu(h_{t-1})$.

Suppose that $h_t$ is a positive probability history in which $\widehat{y}$ occurs in period t. Taking liberties with notation, let $\rho(\widehat{y} \mid \widehat{\Theta}, \alpha^{-1}(h_{t-1}))$ denote the probability of signal $\widehat{y}$ conditional on the unfriendly types. From Bayes rule

$$\mu(h_t)[\widehat{\Theta}] = \frac{\rho(\widehat{y} \mid \widehat{\Theta}^1, \alpha^{-1}(h_{t-1})) \mu(h_{t-1})[\widehat{\Theta}]}{\rho(\widehat{y} \mid \overline{\alpha}^{-1}(h_{t-1}), \alpha^{-1}(h_{t-1}))}$$

Since $\widehat{y}$ has positive probability at time $t$ conditional on $h_{t-1}$, it must be that $\alpha^{-1}(h_{t-1}) \in B(\overline{\alpha}^1(h_{t-1}))$ has positive probability of entry. It follows that $\overline{\alpha}^1(h_{t-1})$ puts weight less than $\widehat{\alpha}$ on $\widehat{A}^1$, and that

$$\frac{\rho(\widehat{y} \mid \widehat{\Theta}, \alpha^{-1}(h_{t-1}))}{\rho(\widehat{y} \mid a^1, \alpha^{-1}(h_{t-1}))} \ge r.$$

Consequently

$$\mu(h_t)[\widehat{\Theta}] \ge \frac{\mu(h_{t-1})[\widehat{\Theta}]}{\widehat{\alpha} + (1 - \widehat{\alpha})\dfrac{1}{r}}.$$

It follows that if signals in $\widehat{Y}$ occur $k$ times, then

$$\mu(h_t)[\widehat{\Theta}] \ge \left( \frac{1}{\widehat{\alpha} + (1 - \widehat{\alpha})\dfrac{1}{r}} \right)^k \mu(0)[\widehat{\Theta}]$$

Hence if

$$k \ge -\frac{\log(\widehat{\alpha})}{\log\left(\widehat{\alpha} + (1 - \widehat{\alpha})\dfrac{1}{r}\right)} - \log\big(\mu(0)[\widehat{\Theta}]\big)$$

then $\mu(h_t)[\widehat{\Theta}] \ge \widehat{\alpha}$, so in all subsequent periods the signal must be an exit signal. This proves the first assertion.

To prove the second part, apply Lemma 1 inductively to conclude that $\alpha_0^1(h_t) \geq (\underline{\alpha}/2)f^1$. Because enforceability is a property of the support of an action, $f^1$ must be enforceable as well as friendly. By assumption every enforceable friendly action is vulnerable to temptation, so there is at least a $\underline{\rho\alpha}/2$ chance of each bad signal $\widehat{y} \in \widehat{Y}$, so when $\widehat{y}$ occurs, from Bayes rule we have $\mu(h_t)[\theta] \leq (1/\underline{\rho\alpha})\mu(h_{t-1})[\theta]$ for all $\theta$. Hence, if $\mu(h_t)[\Theta(F^1)] \geq \underline{\alpha}/2$ then it must be that $\widehat{Y}$ has occurred at least $k$ times, with $(1/\underline{\rho\alpha})^k\mu(0)[\Theta(F^1)] \geq \underline{\alpha}/2$. However, for $\theta(a^1) \notin \widehat{\Theta}$, $\widehat{a}^1 \in \widehat{A}^1$, if $\widehat{y}$ occurs, then

$$\frac{\mu(h_t)[\theta(\widehat{a}^1)]}{\mu(h_t)[\theta(a^1)]} \geq r.$$

Consequently, at $k$, it must be that

$$\mu(h_t)[\widehat{\Theta}] \geq r^k\mu(h_t)[\theta(F^1)] \geq r^k\underline{\alpha}/2.$$

However, under the hypothesis of the Lemma, it can be verified that this implies $\mu(h_t)[\widehat{\Theta}] \geq \widehat{\alpha}$. In other words, as soon as $\mu(h_t)[\Theta(F^1)] \geq \underline{\alpha}/2$, the equilibrium play of the short-run players is concealing for the remainder of the game.

$$\boxtimes$$

**Lemma 3:** *In a participation game if $\alpha^{-1}(h_t) \in \text{conhull}\, E^{-1}$ or $\alpha^{-1}(h_t) \notin \text{conhull}\, E^{-1}$ and $\alpha_0^1(h_t) \geq \gamma f^1$ for some $\gamma > 0$ and vulnerable friendly action $f^1$ of temptation size $\underline{\rho}, \widetilde{\rho}$ then*

$$v^1(h_t) \leq \max_{y \in Y(h_t)}(1 - \delta)\overline{u}^1(y, \underline{\rho}) + \overline{\delta}(y, \widetilde{\rho})v^1(h_t, y).$$

*Proof:* If $\alpha^{-1}(h_t) \in \text{conhull}\, E^{-1}$ then $Y(h_t) \subseteq Y^e$, and the bound follows directly from the definition of $\overline{u}^1$. So consider $\alpha^{-1}(h_t) \notin \text{conhull}\, E^{-1}$. Note that $f^1$ must be enforceable, since otherwise the rational type could replace $f^1$ and get a strictly higher utility. Since every enforceable friendly $f^1$ is vulnerable to temptation, $\rho(\widehat{Y} \mid \overline{\alpha}^1(h_t), \alpha^{-1}(h_t)) \geq \underline{\alpha}\rho(\widehat{Y} \mid f^1, \alpha^{-1}(h_t)) > 0$ so that $\widehat{Y} \subseteq Y(h_t)$. Notice that $v^1(h_t, h_t^1)$ must be independent of the private history $h_t^1$, and that the

rational long-run player must be indifferent between the actions in the support of $\alpha^1(h_t)$. In particular, playing $f^1$ must yield the expected present value $v^1(h_t)$.

Consider then the long-run player switching from playing $f^1$ to a $b^1$ given in the definition of temptation. We now need to calculate the long-run player payoff separately as a function of whether the short-run players exit or not. Observe that $\alpha^{-1}(h_t)$ induces a probability distribution over $A^{-1}$. The probability distribution can be written as a convex combination of two component distributions, namely $\alpha_{-e}^{-1}$, which has support entirely in $A^{-1} - E^{-1}$, and $\alpha_e^{-1}$, which has support entirely in $E^{-1}$. Then $\lambda \alpha_{-e}^{-1} + (1 - \lambda)\alpha_e^{-1}$ induces the same distribution over $A^{-1}$ as $\alpha^{-1}(h_t)$, where $\alpha^{-1}(h_t) \notin \text{conhull } E^{-1}$, $\lambda > 0$. Notice that in general, $\alpha_{-e}^{-1}$ and $\alpha_e^{-1}$ do not have representation as mixed strategies, as they can correspond to correlated strategies for the short-run players. However, we may still write $u^1(\alpha^1, \alpha_e^{-1}), v^1(h_t, \alpha^1, \alpha_e^{-1}), \rho(\cdot \mid \alpha^1, \alpha_e^{-1})$ and so forth for the expected values of $u^1(\alpha^1, a^{-1}), v^1(h_t, \alpha^1, a^{-1}), \rho(\cdot \mid \alpha^1, a^{-1})$ with respect to the weights $\alpha_e^{-1}$, and similarly for $\alpha_{-e}^{-1}$. For example,

$$u^1(\alpha^1, \alpha_e^{-1}) \equiv \sum_{a^{-1}} u^1(\alpha^1, a^{-1})\alpha_e^{-1}[a^{-1}]$$

With this in mind, we may decompose

$$v^1(h_t) = v^1(h_t, f^1, \alpha^{-1}(h_t))$$
$$\geq v^1(h_t, b^1, \alpha^{-1}(h_t)) = \lambda v^1(h_t, b^1, \alpha_e^{-1}) + (1 - \lambda)v^1(h_t, b^1, \alpha_{-e}^{-1})$$

Since $u^1(b^1, \alpha_e^{-1}) \geq u^1(f^1, \alpha_e^{-1})$ by definition of a temptation, it must be that $v^1(h_t, f^1, \alpha_{-e}^{-1}) \geq v^1(h_t, b^1, \alpha_{-e}^{-1})$. Notice that in the current period, $b_1$ loses at most $U^1$. Let $v^1(\alpha^1, \sim \hat{Y})$ denote the continuation expected present value, conditional on $\alpha^1, a_{-e}^{-1}$ and a signal not in $\hat{Y} \cup Y^e$. The overall next period present value of for $f^1$ is

$$\sum_{\hat{y} \in \hat{Y}} \rho(\hat{y} \mid f^1, \alpha_{-e}^{-1})v^1(h_t, \hat{y}) + \rho(\sim \hat{Y} \mid f^1, \alpha_{-e}^{-1})v^1(f^1, \sim \hat{Y}) \ .$$

By switching, the next period present value is

$$\sum_{\hat{y}\in\hat{Y}} \rho(\hat{y} \mid b^1, \alpha_{\sim e}^{-1})v^1(h_t, \hat{y}) + \rho(\sim \hat{Y} \mid b^1, \alpha_{\sim e}^{-1})v^1(b^1, \sim \hat{Y}).$$

Since $f^1$ is in fact optimal, we conclude that

$$(1-\delta)U^1 + \delta\sum_{\hat{y}\in\hat{Y}} \rho(\hat{y} \mid f^1, a_{\sim e}^{-1})v^1(h_t, \hat{y}) + \rho(\sim \hat{Y} \mid f^1, \alpha_{\sim e}^{-1})v^1(f^1, \sim \hat{Y})$$

$$-\delta\sum_{\hat{y}\in\hat{Y}} \rho(\hat{y} \mid b^1, \alpha_{\sim e}^{-1})v^1(h_t, \hat{y}) + \rho(\sim \hat{Y} \mid b^1, \alpha_{\sim e}^{-1})v^1(b^1, \sim \hat{Y}) \geq 0$$

or

$$(1-\delta)U^1 + \delta\sum_{\hat{y}\in\hat{Y}}\left[\rho(\hat{y} \mid f^1, \alpha_{\sim e}^{-1}) - \rho(\hat{y} \mid b^1, \alpha_{\sim e}^{-1})\right]v^1(h_t, \hat{y}) \geq$$

$$\delta\left[\sum_{y\in\sim\hat{Y}}\left[\rho(y \mid b^1, \alpha_{\sim e}^{-1}) - \rho(y \mid f^1, \alpha_{\sim e}^{-1})\right]v^1(h_t, y)\right]$$

But $b^1$ reduces the probability of every bad signal by at least $\underline{\rho}$, and the continuation payoff $v^1(h_t, y)$ must be at least the minmax value, which we have set to 0. This implies for

$$v^1(h_t, \hat{y}) = \max_{y\in\hat{Y}} v^1(h_t, y)$$

it must be that

$$(1-\delta)U^1 + \delta v^1(h_t, \hat{y}) \geq$$

$$\delta\left[\rho(\sim \hat{Y} \mid b^1, \alpha_{\sim e}^{-1})v^1(b^1, \sim \hat{Y}) - \rho(\sim \hat{Y} \mid f^1, \alpha_{\sim e}^{-1})v^1(f^1, \sim \hat{Y})\right]$$

Continuing to use the fact that $v^1(h_t, y) \geq 0$, part 2 of the definition of a temptation implies

$$\delta\sum_{y\in\sim\hat{Y}}\left[\rho(y \mid b^1, \alpha_{\sim e}^{-1}) - \rho(y \mid f^1, \alpha_{\sim e}^{-1})\right]v^1(h_t, y)$$

$$\geq \delta\tilde{\rho}v^1(f^1, \sim \hat{Y})$$

from which

$$v_{\sim e}^1 \equiv (1-\delta)u^1(f^1, \alpha_{\sim e}^{-1}) + \delta\sum_{\hat{y}\in\hat{Y}} \rho(\hat{y} \mid f^1, a_{\sim e}^{-1})v^1(h_t, \hat{y})$$

$$+\rho(\sim \hat{Y} \mid f^1, a_{\sim e}^{-1})v^1(f^1, \sim \hat{Y})$$

$$\leq (1-\delta)U^1 + \frac{(1-\delta)U^1}{\tilde{\rho}} + \frac{\delta}{\tilde{\rho}}v^1(h_t, \hat{y})$$

$$= \max_{y\in\hat{Y}} (1-\delta)\bar{u}^1(y, \underline{\rho}) + \frac{\delta}{\tilde{\rho}}v^1(h_t, y)$$

If $\lambda = 1$ this gives the desired result. If not,

$$v^1(h_t) = \lambda v^1(h_t, f^1, \alpha_e^{-1}) + (1-\lambda)v_{\sim e}^1$$

$$\leq \max\{v^1(h_t, f^1, \alpha_e^{-1}), v_{\sim e}^1\}$$

Since

$$v^1(h_t, f^1, \alpha_e^{-1}) \leq \max_{y \in Y(h_t)}(1-\delta)\bar{u}^1(y, \underline{\rho}) + \bar{\delta}(y, \tilde{\rho})v^1(h_t, y)$$

the result follows

☑

**Lemma 4:** *In a participation game, if $\alpha^{-1}(h_t) \in \text{conhull}\, E^{-1}$ or $\alpha^{-1}(h_t) \notin \text{conhull}\, E^{-1}$ and $\alpha_0^1(h_t) \geq \gamma f^1$ for some $\gamma > 0$ and friendly action $f^1$ that is vulnerable to a strong temptation size $\underline{\rho}$, then*

$$v^1(h_t) \leq \max_{y \in Y(h_t)}(1-\delta)\bar{u}^1(y, \underline{\rho}) + \delta v^1(h_t, y).$$

*Proof:* Since the proof of Lemma 4 is very similar to that of lemma 3, we will only discuss the necessary changes. As in the proof of Lemma 3, we consider the long-run player switching from playing $f^1$ to a $b^1$ given in the definition of what is now a strong temptation. Using the notation of the proof of lemma 3, we compute that the overall next period present value of for $f^1$ is

$$\sum_{\hat{y} \in \hat{Y}} \rho(\hat{y} \mid f^1, a_{\sim e}^{-1})v^1(h_t, \hat{y}) + \rho(\sim \hat{Y} \mid f^1, a_{\sim e}^{-1})v^1(\sim \hat{Y}).$$

Because of the proportionality of the probabilities for a strong temptation, the corresponding value for $b^1$ is

$$\sum_{\hat{y} \in \hat{Y}} \rho(\hat{y} \mid b^1, a_{\sim e}^{-1})v^1(h_t, \hat{y}) + \rho(\sim \hat{Y} \mid b^1, a_{\sim e}^{-1})v^1(\sim \hat{Y}).$$

Since $f^1$ is in fact optimal, we conclude that

$$(1-\delta)U^1 + \delta\sum_{\hat{y} \in \hat{Y}} \rho(\hat{y} \mid f^1, a_{\sim e}^{-1})v^1(h_t, \hat{y}) + \rho(\sim \hat{Y} \mid f^1, a_{\sim e}^{-1})v^1(\sim \hat{Y})$$

$$-\delta\sum_{\hat{y} \in \hat{Y}} \rho(\hat{y} \mid b^1, a_{\sim e}^{-1})v^1(h_t, \hat{y}) + \rho(\sim \hat{Y} \mid b^1, a_{\sim e}^{-1})v^1(\sim \hat{Y}) \geq 0$$

or

$$(1 - \delta)U^1 + \delta \sum_{\hat{y} \in \hat{Y}} \left[ \rho(\hat{y} \mid f^1, a_{\sim e}^{-1}) - \rho(\hat{y} \mid b^1, a_{\sim e}^{-1}) \right] v^1(h_t, \hat{y}) \geq$$

$$\delta \left[ \rho(\sim \hat{Y} \mid b^1, a_{\sim e}^{-1}) - \rho(\sim \hat{Y} \mid f^1, a_{\sim e}^{-1}) \right] v^1(\sim \hat{Y})$$

Set

$$v^1(h_t, \hat{Y}) = \max_{y \in \hat{Y}} v^1(h_t, y).$$

From the fact that $b^1$ reduces the probability of every bad signal by a positive amount

$$(1 - \delta)U^1 + \delta \left[ \rho(\hat{Y} \mid b^1, a_{\sim e}^{-1}) - \rho(\hat{Y} \mid f^1, a_{\sim e}^{-1}) \right] v^1(h_t, \hat{Y}) \geq$$

$$\delta \left[ \rho(\hat{Y} \mid b^1, a_{\sim e}^{-1}) - \rho(\hat{Y} \mid f^1, a_{\sim e}^{-1}) \right] v^1(\sim \hat{Y})$$

Observe that since $b^1$ reduced the probability of each bad signal by at least $\underline{\rho}$, $\left[ \rho(\hat{Y} \mid b^1, a_{\sim e}^{-1}) - \rho(\hat{Y} \mid f^1, a_{\sim e}^{-1}) \right] \geq \# \hat{Y} \underline{\rho}$. Since $U^1 > 0$ it must be that

$$(1 - \delta)U^1 + \delta \, \# \hat{Y} \underline{\rho} v^1(h_t, \hat{y}) \geq \delta \, \# \hat{Y} \underline{\rho} v^1(\sim \hat{Y})$$

from which

$$v_{\sim e}^1 \equiv (1 - \delta)u^1(f^1, \alpha_{\sim e}^{-1})$$

$$+ \delta \sum_{\hat{y} \in \hat{Y}} \rho(\hat{y} \mid f^1, a_{\sim e}^{-1}) v^1(h_t, \hat{y}) + \rho(\hat{Y} \mid f^1, a_{\sim e}^{-1}) v^1(\sim \hat{Y})$$

$$\leq (1 - \delta)U^1 + \frac{(1 - \delta)U^1}{\# \hat{Y} \underline{\rho}} + \delta v^1(h_t, \hat{y}) =$$

$$\max_{y \in \hat{Y}} (1 - \delta)\bar{u}^1(y, \underline{\rho}) + \delta v^1(h_t, y)$$

The final steps exactly parallel those of the proof of Lemma 3.

☑

# References

Ely, J. and J. Valimaki [2000] "Bad Reputation," mimeo.

Fudenberg, D. and D. Kreps [1987] "Reputation and Simultaneous Opponents" *Review of Economic Studies*, 54: 541-568

Fudenberg, D., D. Kreps, and E. Maskin [1990] "Repeated Games with Long-run and Short-run Players," *Review of Economic Studies*, 57, 555-573.

Fudenberg, D. and D. K. Levine [1994] "Efficiency and Observability in Games with Long-Run and Short-Run Players," *Journal of Economic Theory,* 62 , 103-135

Fudenberg, D. and D. K. Levine [1992] "Maintaining a Reputation when Strategies are Imperfectly Observed," *Review of Economic Studies,* 59: 561-579.

Fudenberg, D. and D. K. Levine [1989] "Reputation and Equilibrium Selection in Games with a Single Long-Run Player" *Econometrica*, 57: 759-778.

Fudenberg, D., E. Maskin, and D.K. Levine [1994] "The Folk Theorem in Repeated Games with Imperfect Public Information," *Econometrica* 62, 997-1039.

Kreps, D. and R. Wilson [1982] "Reputation and Imperfect Information," *Journal of Economic Theory,* 27:253-279.

Mailath, G. and L. Samuelson [1998] "Your Reputation is Who You're Not, Not Who You'd like to Be," mimeo.

Milgrom, P. and J. Roberts [1982] "Predation, Reputation, and Entry Deterence," *Journal of Economic Theory,* 27:280-213.

Sorin, S. [1999] Merging, Reputation, and Repeated Games with Incomplete Information," *Games and Economic Behavior* 29, 274-308