

Order Stability in Supply Chains: Coordination Risk and the Role of Coordination Stock

Rachel Croson^{*}, Karen Donohue^{**}, Elena Katok⁺, and John Sterman⁺⁺

^{*}School of Economics, University of Texas at Dallas, Richardson, TX, 75080-3021,
crosnr@utdallas.edu;

^{**}The Carlson School, University of Minnesota, Minneapolis, MN 55455-9940, donoh008@umn.edu;

⁺Jindal School of Management, University of Texas at Dallas, Richardson, TX , 75080-3021,
ekatok@utdallas.edu;

⁺⁺Sloan School of Management, Massachusetts Institute of Technology, Cambridge, MA 02142,
jsterman@mit.edu

Abstract

The *bullwhip effect* describes the tendency for the variance of orders in supply chains to increase as one moves upstream from consumer demand. Previous research attributes this phenomenon to both operational and behavioral causes. We report on a set of laboratory experiments with a serial supply chain, using the Beer Distribution Game. The experimental conditions eliminate all operational causes of the bullwhip effect. Nevertheless, we find that the bullwhip effect persists in this setting and offer one possible explanation based on coordination risk. Coordination risk exists when individuals' decisions contribute to a collective outcome and the decision rules followed by each individual are not known with certainty, e.g., where managers cannot be sure how their supply chain partners will behave. We conjecture that the existence of coordination risk may contribute to bullwhip behavior. We test this conjecture by controlling for environmental factors that lead to coordination risk and find these controls lead to significant reduction in order oscillations and amplification. Next, we investigate a managerial intervention to reduce the bullwhip effect, inspired by our conjecture that coordination risk contributes to bullwhip behavior. While the intervention, holding additional on-hand inventory, does not change the existence of coordination risk, we find that it reduces order oscillation and amplification by providing a buffer against the endogenous risk of coordination failure. We conclude that the magnitude of the bullwhip can be mitigated, but that its behavioral causes appear robust.

Keywords: Bullwhip Effect, Behavioral Operations, Supply Chain Management, Beer Distribution Game

Revised, May 2012

1. Introduction

Despite the undoubted benefits of the lean manufacturing and supply chain revolutions of the past decade, supply chain instability continues to plague many businesses (Ellram 2010). The consequences include excessive inventories, periodic stockouts, poor customer service, and unnecessary capital investment. These outcomes have recently been documented in many manufacturing and wholesale sectors of the economy (e.g., Aeppel 2010, Dooley et al. 2010).

Supply chain instability is often attributed, in part, to the *bullwhip effect*. The bullwhip effect refers to the tendency for order variability to increase within a supply chain as orders move upstream from customer sales to production. Lee et al. (1997) and Sterman (2000) cite evidence of this phenomenon in a wide range of industries. Cachon, Randall, and Schmidt (2007) find that the phenomenon is most prevalent in wholesale industries and industries with low seasonality. Chen and Lee (2011) provide a detailed analysis of how the level of time aggregation used in measuring the bullwhip effect can influence whether the effect is detected or not. They also provide insight into what level of aggregation is most appropriate for capturing the connection between bullwhip behavior and supply chain cost.

Previous research attributes the bullwhip effect to both operational and behavioral causes. Operational causes are structural characteristics that lead rational agents to amplify demand variation. Examples include order batching, gaming due to shortages, price fluctuations caused by promotions, and demand signaling due to forecast uncertainty (Lee et al. 1997). These causes have been documented in practice, and techniques to eliminate them are now an important part of the tool kit for supply chain management (Lee et al. 2004). Behavioral causes, in contrast, emphasize the bounded rationality of decision makers, particularly the failure to adequately account for feedback effects and time delays. Specifically, people tend to place orders based on the gap between their target inventory level and current on-hand stock, while giving insufficient weight to the supply line of unfilled orders (the stock of orders placed but not yet received). Prior work demonstrates that the tendency to underweight the supply line is sufficient to cause the bullwhip effect in both experimental and real supply chains (e.g., Sterman 1989a, Sterman 2000).

The goal of this paper is to extend our understanding of the behavioral causes of the bullwhip effect by examining how individuals perform when all operational causes of the bullwhip effect are eliminated. In the absence of operational causes, any remaining supply chain instability and bullwhip must arise from behavioral causes. We test this hypothesis in a baseline experiment, which follows the general protocol of the beer distribution game except that demand uncertainty is eliminated by using a constant demand rate that is publicly announced to all participants.¹ The experimental setup significantly simplifies the

¹ The standard protocol for the beer game eliminates the need for order batching because there are no fixed ordering costs or quantity discounts. Each player orders from and ships to only one other party and production capacity is infinite, eliminating the incentive for gaming in response to shortages. Prices are fixed and customer demand is exogenous, eliminating the possibility of promotions and forward buying. In our design, we also eliminate demand signaling since demand is constant and known to all participants.

decisions participants must make compared to real supply chains as well as prior experimental research, by effectively eliminating all previously cited operational causes of the bullwhip effect. Also, because demand is fixed at a known and constant rate, there is no need for safety stock, eliminating possible cognitive errors in computing the appropriate safety stock level that might occur with a more complex demand distribution. Additionally, the system is started in equilibrium, at the optimal on-hand inventory level of zero, so there is no need for participants to calculate the optimal transient path to reach an equilibrium level of inventory.

Nevertheless, we find that order oscillation and amplification persist even in this simplified environment. Estimation of participants' decision rules shows that the vast majority underweight the supply line of unfilled orders, consistent with prior studies. Post-play questionnaires and participant debriefing suggest that many players were unable to predict how their teammates would behave.

The results suggest that the existence of *coordination risk* may partly account for the observed behavior. Coordination risk exists when individuals in a group environment make independent decisions that contribute to the collective outcome and the decision rules followed by each individual are not known with certainty. Such uncertainty arises when decision-makers have limited knowledge of, or trust in, their partners' motives or cognitive abilities. This uncertainty may cause decision makers to deviate from equilibrium to hedge against the possibility that their partners will not behave optimally. In real supply chains, inventory managers often have limited understanding of each other's decision rules or local incentives. Such uncertainty is particularly likely when the optimal decision rule is not commonly known.

The next two experiments test the effect of eliminating environmental factors that may lead to coordination risk by first creating common knowledge about the optimal ordering policy (experiment 2) and then ensuring that all supply chain partners of the focal human subject follow the optimal policy by automating their decisions (experiment 3). We find that order oscillation and amplification are significantly reduced in these experiments compared with the baseline. However, the large majority of participants continue to underweight the supply line.

Lastly, in experiment 4 we test a managerial intervention to test our conjecture that coordination risk contributes to bullwhip behavior. While the intervention, providing participants with additional on-hand inventory, does not eliminate coordination risk, we find that it reduces order oscillation and amplification by providing a buffer against the endogenous risk of coordination failure. The results suggest a new purpose for inventory, which we term *coordination stock*. Coordination stock differs from conventional safety stock held to buffer against exogenous uncertainty in demand or deliveries. In contrast, coordination stock is inventory held to buffer against strategic uncertainty about the actions of other supply chain members. We find that a modest level of coordination stock significantly lowers order variability and costs, suggesting a rationale for what may appear to be excess inventory in decentralized

supply chains.

Overall, our results support the role of behavioral causes in the bullwhip effect. Coordination risk and supply line underweighting are likely both present and contribute to deviation from optimal behavior. Directly reducing the strategic uncertainty that contributes to coordination risk significantly improves performance, although continued supply line underweighting means that oscillation and amplification still exist. Holding coordination stock does not reduce coordination risk directly, but instead alleviates its effect, and thus also has a positive effect on performance. That said, supply line underweighting is highly robust—it does not vary significantly by supply chain role and is only slightly moderated by the experimental treatments.

We continue in the next section with a brief overview of prior literature and its connection to our main hypotheses. Section 3 provides details of the experimental protocol, followed in section 4 by a description of the data generated by the experiments and the associated analysis. Section 5 provides a general discussion of the results and possible psychological processes behind the behavior we observed. Section 6 concludes with potential managerial implications and directions for future research.

2. Theory Development and Hypotheses

Our research fits within the growing field of behavioral operations management. This field aims to identify, analyze, and fill gaps in understanding between operations theory and actual human behavior. Such work requires adapting tools and theories from social science disciplines, such as experimental economics, psychology, system dynamics, and organizational behavior. The behavioral operations management literature relating to our present study examines inventory ordering decisions in multi-player settings, including supply chains. In this stream of research, summarized in Croson and Donohue (2002), researchers have focused on estimating individual decision rules for ordering and comparing them with optimal rules, e.g., Sterman (1989a, 1989b). Much of this research has been conducted in the context of the beer distribution game. The game is well-suited to examine behavioral issues since it is simple enough for people to learn quickly while retaining key features of real supply chains.

The major difference between our baseline setting (experiment 1) and prior work using the beer distribution game is the treatment of customer demand. In previous work, customer demand both varied across time periods (i.e., was either stochastic or non-stationary), and was unknown to participants. For example, Sterman (1989a) and Steckel et al. (2004) use unknown, non-stationary but deterministic demand functions (a step function and S-curve, respectively). Croson and Donohue (2003, 2006) use a known, stationary uniform distribution, but agents (other than the retailer) do not know the realizations of demand. Due to the order fulfillment delay, agents must forecast future demand. It is conceivable that the bullwhip in these settings results, in part, from demand uncertainty and resulting forecasting errors (Chen

et al. 2000).

We eliminate demand variability and forecasting as potential causes of the bullwhip by using a constant demand of 4 cases per period and publicly informing participants of this fact before the game begins. When demand is constant and publicly known, there is no need for safety stock. Further, unlike prior studies, we initialize the system in equilibrium so there is no incentive for participants to adjust the current stock level.

By adding these controls, our baseline experiment is designed to test whether the bullwhip effect still exists when all operational causes and operational sources of uncertainty are removed. The only uncertainty facing participants lies in the decision processes other supply chain team members use to make their order decisions. However, it is trivial to show that costs are minimized when each decision maker simply passes through the order they receive. Under this policy, the system remains in equilibrium and there is no order oscillation or amplification (i.e., no bullwhip). Hence, theory suggests:

Hypothesis 1: *The bullwhip effect (order oscillation and amplification) will not occur when demand is known and constant and the system begins in equilibrium.*

This (null) hypothesis implicitly assumes that either coordination risk does not exist or, if it does exist, it does not affect ordering behavior.

The importance of coordination in collective decision settings is well studied in the context of coordination games (e.g., Van Huyck, Battalio and Beil 1990, 1991, 1993, Cachon and Camerer 1996). However, that literature focuses primarily on games with multiple equilibria, in which the failure to coordinate causes players to converge to the risk-dominant, instead of the payoff-dominant equilibrium. A common goal in that literature, which we share, is to identify and test mechanisms designed to help players coordinate.

Unlike the traditional coordination games, our supply chain setting does not have multiple equilibria. The equilibrium is unique, but time delays, accumulations, and feedbacks in the supply chain cause the equilibrium to be potentially unstable. Failure to coordinate may induce instability, oscillation, and other disequilibrium dynamics that complicate or prevent a return to equilibrium. Thus our supply chain setting uncovers new aspects of the coordination problem. In our setting, coordination risk may arise from two factors: (1) a lack of common knowledge about the decision rule that keep the system in equilibrium, and (2) a lack of trust that the other players will implement that rule even if it is commonly known.

The concept of common knowledge originated in game theory (Aumann 1976) and is considered a necessary component of equilibrium in games. Vanderschraaf and Sillari (2005) provide a comprehensive description of the concept and its history in academic research. Agents often act on the basis of incomplete information; decision makers must always determine what they know and what they do not

know before making a decision. When the outcome involves more than one agent, each decision maker must also consider what others may know. Common knowledge goes beyond mutual knowledge. Under mutual knowledge, each agent privately knows the same information. In our case, this information would be the optimal ordering policy. Mutual knowledge is typically implemented in the lab by privately informing all participants about relevant information. Common knowledge goes a step further by requiring not only that agents know the information but that they also know that all the other players also know this information. In addition, players know that other players also know that everyone knows this information and so on *ad infinitum*. Common knowledge is typically implemented in the laboratory by publicly providing the relevant information to participants. After the public announcement, everyone knows that everyone else has been informed, etc.

Game-theoretic models and experimental studies show that failures of common knowledge can lead to coordination failures and inefficient outcomes in a variety of settings (e.g., Geanakoplos 1992, Lei, Noussair, and Plott 2001, Nagel 1995). We contribute to this literature by testing the effect of creating common knowledge on reducing bullwhip behavior, which is a natural outcome of reducing coordination risk. If the bullwhip is caused or triggered by coordination risk, then providing common knowledge may reduce both order variability and amplification (compared to the baseline experiment). This leads to our first hypothesis concerning coordination risk.

Hypothesis 2: *The provision of common information through the announcement of the optimal policy will reduce the bullwhip effect (order oscillation and amplification).*

Also, since the description of the optimal policy explicitly directs participants to consider both on-hand and on-order inventory, theory suggests that decision-making at the individual level may improve by leading agents to fully account for the supply line. Hence,

Hypothesis 3: *The provision of common information through the announcement of the optimal policy will reduce supply line underweighting.*

The second element of coordination risk is *trust* that other agents will actually carry out the optimal ordering policy. While others have studied the effect of a lack of common knowledge on experimental outcomes, they have used settings sufficiently simple that trust is typically not an issue (see Ochs 1995 for a review). A large literature in economics and organizational behavior examines the effect of trust in others' benign intentions (see, e.g. Rousseau et al., 1998). Here, however, we focus not on intentions but on trust in the ability of other players to implement optimal decisions correctly. The closest previous research investigating this type of trust is Schwierien and Sutter (2008), who compare individuals' trust in others' good intentions toward them with trust in others' abilities.

If our explanation for the cause of order oscillation and amplification in the baseline is correct,

eliminating the possibility of coordination failure should eliminate (or at least significantly reduce) order oscillation and amplification. *Overall* supply chain performance will improve by definition, since the automated agents play optimally. The relevant comparison is thus how well the human players perform here compared with the retailers from previous experiments, who face the same incoming order stream. With the guarantee that all others will play optimally, the performance of the human players should improve significantly, and we should see less order oscillation.

Hypothesis 4: *Eliminating coordination risk will decrease order oscillation for the human players, relative to the retailers in the (a) baseline and (b) common knowledge experiments.*

Finally, since the supply line is now more predictable, it should be easier to track. Also, freed of the need to devote scarce cognitive resources to anticipating how their teammates will behave, individuals may be more likely to use the optimal policy and fully account for the supply line. Hence

Hypothesis 5: *Eliminating coordination risk will reduce supply line underweighting.*

Following a common goal of prior coordination research, it is valuable to consider what mechanisms could be used to mitigate the effects of coordination risk when the source of the risk itself cannot be removed. One possible policy intervention is to hold excess inventory to buffer against the interaction effects caused by the possible sub-optimal decisions made by other supply chain members. Such a stabilizing benefit has been suggested by others (e.g., Cohen and Baganha 1998), but never tested in an experimental setting. Previous research has identified operational factors affecting the optimal level of safety stock, including the length of the replenishment lead-time, level of demand variability, and the relative cost of holding inventory versus carrying a backlog. These are the common inputs used to manage inventory levels in most continuous or periodic review inventory systems (Zipkin 2000). Our research contributes to this literature by testing a new role for excess inventory—to buffer against strategic uncertainty about other players' decisions (as opposed to exogenous uncertainty about market demand, production lags or other operational factors).

Additional on-hand inventory may improve performance for two reasons. First, initializing the system with positive on-hand inventory reduces the need for players who seek a buffer stock to increase orders above equilibrium, reducing the incidence and magnitude of deviations that may trigger the instability caused by supply line underweighting. Second, initial on-hand inventory reduces the likelihood that the system will enter the unstable backlog regime, moderating the oscillations generated by players who underweight the supply line. Hence,

Hypothesis 6: *Excess initial inventory will decrease the bullwhip effect (order oscillation and amplification).*

On the other hand, because excess inventory necessarily initializes the system out of inventory equilibrium, the behavior of the other members of the supply chain may actually be more difficult to predict *a priori*. Rational agents would attempt to reduce their inventory levels, but if they or their supply chain partners underweight the supply line, the result could be instability and oscillations, perhaps even larger than those observed in the baseline experiment. Such behavior would refute Hypothesis 6.

3. Experimental Design

Our experiments follow the standard protocol for the Beer Distribution Game. The experiments were run using a computer network with an interface written in Visual Basic. Complete instructions are available in an on-line supplement.² All sessions were conducted at a large public university in the Northeastern United States in the spring semester of 2003.

The Beer Distribution Game consists of four agents (retailer, wholesaler, distributor and factory), each making ordering decisions for one link in a serial supply chain with exogenous final customer demand (Figure 1). The chain operates as a multi-echelon inventory system in discrete time with demand backlogging, infinite capacity, and both order processing and shipping lags (see Chen 1999 and Clark and Scarf 1960 for classic models of such systems). Each period (corresponding to one week) all players experience the following sequence of events: (1) shipments from the upstream decision-maker are received and placed in inventory, (2) incoming orders are received from the downstream decision-maker and either filled (if inventory is available) or placed in backlog, and (3) a new order is placed and passed to the upstream player. Croson and Donohue (2006) and Sterman (1989a) provide full details.

As in most previous studies, the retailer, wholesaler, and distributor (R, W, D) face a two week lag between placing an order and fulfillment by their immediate supplier, and an additional two week transportation lag between fulfillment by the supplier and delivery (a total lag of 4 weeks, assuming the supplier is in stock). The factory (F) has a one-week lag between placing orders (setting the production schedule) and production starts, and a two-week production lag (a total lag of 3 weeks). As in earlier studies, costs are \$0.50/week for each unit held in inventory, and \$1/week for each unit backlogged.

We eliminated demand variability and forecasting as potential causes of the bullwhip in all four experiments by using a constant demand of 4 cases per week and publicly informing participants of this fact before the game began. Further, the realization of demand each week (4 cases) was displayed on every agent's screen at all times, confirming visually that demand is indeed constant. After the rules were explained, but before play began, participants were given a quiz to confirm they understood the rules and calculations, as well as the fact that customer demand would be constant at 4 cases per week for the entire game. All participants correctly answered these questions, confirming that they understood the rules and

²Instructions, the quiz and other materials can be found at www.personal.psu.edu/exk106/Appendix_Compedium_Instructions.pdf.

operation of the game.

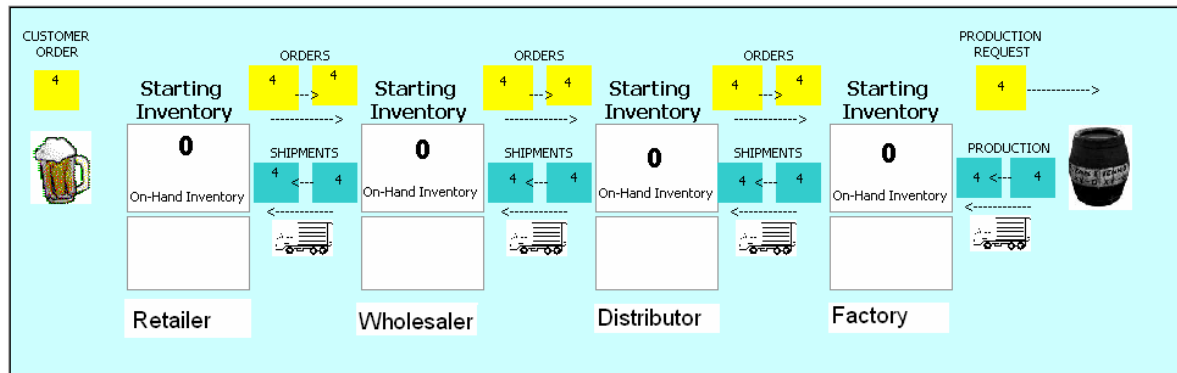


Figure 1: The Beer Distribution Game portrays a serial supply chain.

All experiments began in flow equilibrium with orders and shipments of four cases at each delay step. The initial end-of-period on-hand stock level (after deliveries are received and orders filled) was set to zero for all supply chain levels. No time limit was imposed for making order decisions during the game. All games were run for 48 weeks, although this information was not shared with the participants to avoid horizon effects. After the game, but before receiving payment, participants were asked to complete an on-line questionnaire inviting them to reflect on their experience.

The four experiments varied in the level of common knowledge of the optimal order policy, the automation of supply chain partner decisions, and the initial level of on-hand inventory. Table 1 summarizes these differences. In experiment 2 and 3 we provide common knowledge about the optimal ordering policy, in addition to common knowledge of customer demand (which is provided in all four experiments). This is accomplished by providing the following public explanation of the optimal policy at the beginning of the experiment:

“The total team cost can be minimized if all team members place orders so as to make the total of their on-hand inventory and outstanding orders equal to a pre-specified target level. This target level is 16 for the retailer, wholesaler and distributor, and 12 for the factory. This means that if the total on-hand inventory and outstanding orders is greater than or equal to the target level, the order that minimizes team cost is 0. But if this total inventory is less than the target level, the cost-minimizing order is to order just enough to bring it to the target.”

The explanation was presented in written instructions and explained publicly along with the game rules. The explanation places no restrictions on participants’ actions, but the fact that the explanation is public provides common knowledge: each player is informed of the cost minimizing policy and knows that every other player is also so informed.

In experiment 3 we further control the possible sources of coordination risk by placing individuals in a supply chain with three automated players, each programmed to use the optimal policy. We ran 4 sets of 10 supply chains with each set placing the human decision maker in a different role (i.e., 10

experiments with a human retailer, 10 with a human wholesaler, etc.). Participants are publicly told the optimal, cost minimizing, decision rule as in the previous experiment, and also that all other members of their supply chain are automated and programmed to follow that rule. All other conditions are identical to the previous two experiments.

	Experiment 1 (Baseline)	Experiment 2 (Creating Common Knowledge)	Experiment 3 (Eliminating Coordination Risk)	Experiment 4 (Adding Coordination Stock)
Publicly known and constant demand	Yes	Yes	Yes	Yes
Initial on-hand inventory set at equilibrium level	Yes	Yes	Yes	No (set to 12)
Optimal order rule publicly announced	No	Yes	Yes	No
One's fellow supply chain members guaranteed to follow optimal policy (automated)	No	No	Yes	No
Undergraduate*	90%	94.7%	98.7%	90.7%
Gender – female	56.9%	31.6%	47.4%	39.5%
Major – business	55.2%	63.2%	57.9%	72.1%
Major – science & engineering	15.5%	10.5%	15.8%	9.3%
Major – liberal arts, communications, and education	19.0%	21.1%	23.7%	14.0%
Major – other**	10.3%	5.2%	2.6%	4.6%
N (individuals)	40	40	40	40
T (teams)	10	10	40	10
N and T (minus outliers)	32, 8	40, 10	35, 35	28, 7

* All others are graduate students.

** Other majors include HDFS, Agricultural Science, IST, EMS, Arts and Architecture, and undecided.

Table 1: Conditions and Demographic Breakdown Across Experiments

Note: Although participants were randomly assigned to experiments *ex-ante*, we ran t-tests of proportions to check if there are systematic differences *ex-post* in these demographics by experiment. Out of the 36 comparisons run (each experiment versus each other experiment on each dimension reported above) we find only one significant difference (between the proportion of women in the baseline and common knowledge experiments). When running multiple tests, it is not uncommon to get false positive results (in fact, 1 in 20 tests will show a 5% significant difference by chance). In these situations, a Bonferroni correction is typically used. Once this correction is applied this gender difference is no longer statistically significant

In experiment 4, conditions are identical to the baseline except that all participants begin with 12 units of on-hand inventory rather than the optimal level of zero. Twelve units is the initial inventory traditionally used in the beer game, facilitating comparison to prior work, and likely exceeds the desired level of coordination stock (as suggested by the estimates of S' in Table 2). Our purpose is not to find the optimal level of coordination stock, but rather to test whether directionally there is an improvement when excess inventory is provided.

Table 1 also summarizes participant demographics for each treatment. A total of 160 participants were recruited using an on-line recruitment system. All participants were drawn from the same student pool and randomly assigned to one of the four treatments. Participants were paid a \$5 show-up fee, plus up to \$20 in bonus money depending on their team’s relative performance. Team performance was computed based on total supply chain cost (i.e., the cumulative sum of holding and backorder costs across all four players) using the continuous payment scheme introduced by Croson and Donohue (2006).

We tested the data from each treatment for outliers using the Grubbs procedure (Grubbs 1969). A small number of outliers were found and deleted from the statistical analyses reported in the main text (the Appendix details the treatment of outliers). Eliminating the outliers reduces the between-condition differences, an *a fortiori* procedure that provides a tougher test of our hypotheses. For completeness, we also report results using the entire data set (i.e., including outliers) in footnotes. All our findings hold at similar levels of significance with or without outliers.

4. Experimental Results

4.1 Experiment 1: Baseline with Constant Demand

Figure 2 shows orders and on-hand inventory for each individual in our baseline experiment. The characteristic pattern of oscillation and amplification observed in prior research persists. None of the teams remain close to the equilibrium throughput of 4 cases/week. Excluding outliers, factory orders average 19; the mean of the peak factory orders is 215 and the maximum factory order is 500.

Figure 3 displays the standard deviations of orders for each role across the eight teams in our data analysis (teams 4 and 9 were identified as outliers and excluded). Significant order oscillation clearly exists despite constant demand and common knowledge of that fact. We use a nonparametric sign test to check for order amplification (Seigel 1965, p. 68). Order amplification exists when the standard deviation of the i^{th} stage in the supply chain, σ_i , exceeds that of its immediate customer, σ_{i-1} (where $i \in \{R, W, D, F\}$). If there were no bullwhip, we would observe $\sigma_i > \sigma_{i-1}$ at the chance rate of 50%. The data reveal $\sigma_i > \sigma_{i-1}$ for 79% of the cases, rejecting H1 at $p = 0.0025$.³ The bullwhip effect persists even when customer demand is constant and publicly known to all participants.⁴

4.1.1 Individual Ordering Behavior

To better understand the decision processes of the participants, we specify a decision rule for managing inventory and estimate its parameters for each participant. Comparing the estimated decision weights to the optimal weights offers insight into participant decision-making.

³ Including the outlier teams (4 and 9), $\sigma_i > \sigma_{i-1}$ for 80% of the cases, and H1 is rejected at $p = 0.0005$.

⁴ To test the robustness of our results to the use of the standard deviations of orders as the metric, we also analyzed the root mean squared (RMS) deviation, δ , of orders, O , from the steady-state optimal value of 4, $\delta_{i,j} = (1/n)[\sum_{t \in \{1, 48\}} (O_{i,j,t} - 4)^2]^{1/2}$ for role i in team j . The RMS deviation penalizes participants for all deviations from optimal, and not only for variability around the participant’s mean orders. All statistical results reported here continue to hold using the RMS deviation from optimal instead of the standard deviation.

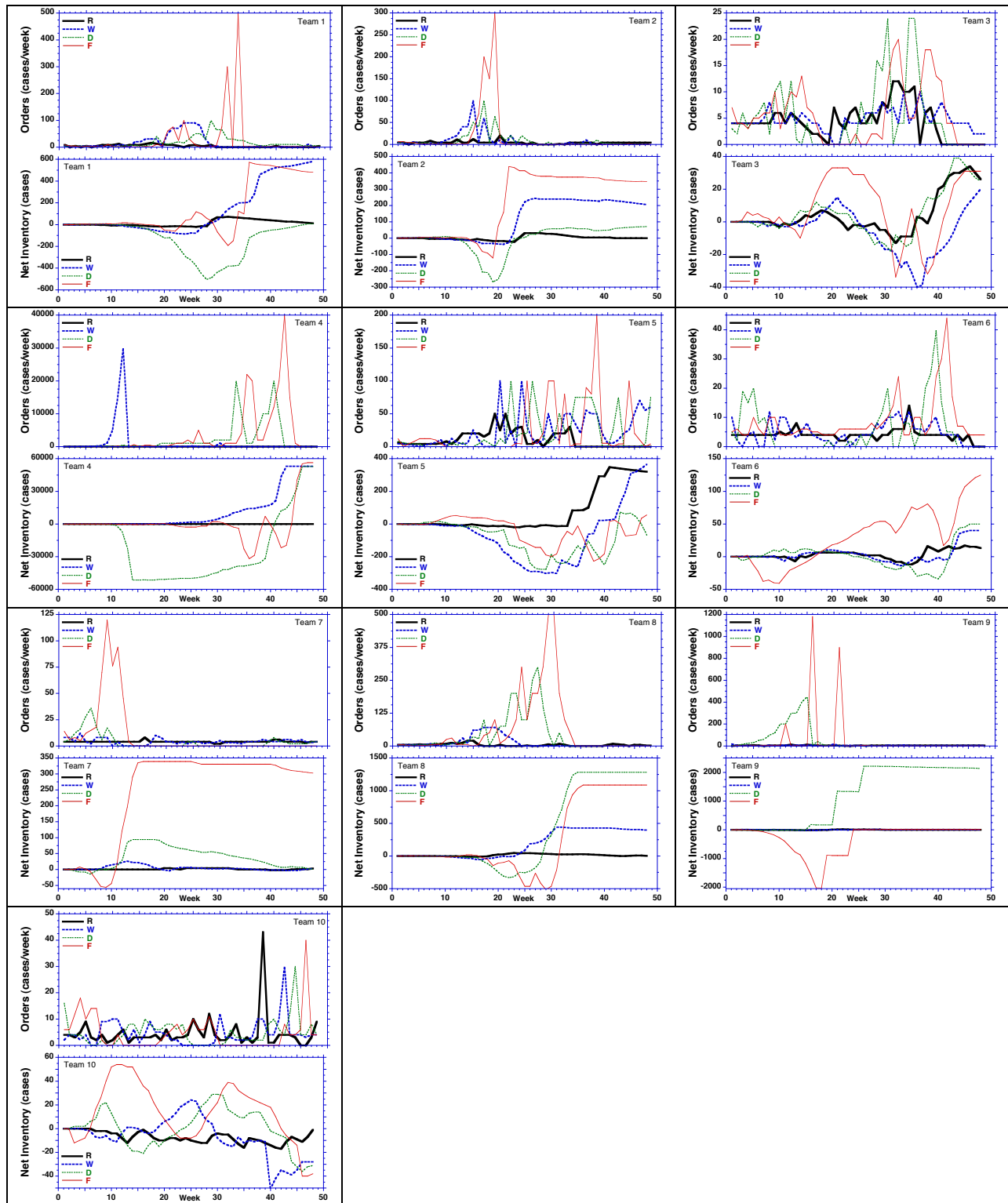


Figure 2. Orders (top) and On-Hand Inventory (bottom) for the 10 teams in Experiment 1 (Baseline). Negative inventory values indicate backlogs. Note: vertical scales differ across teams.

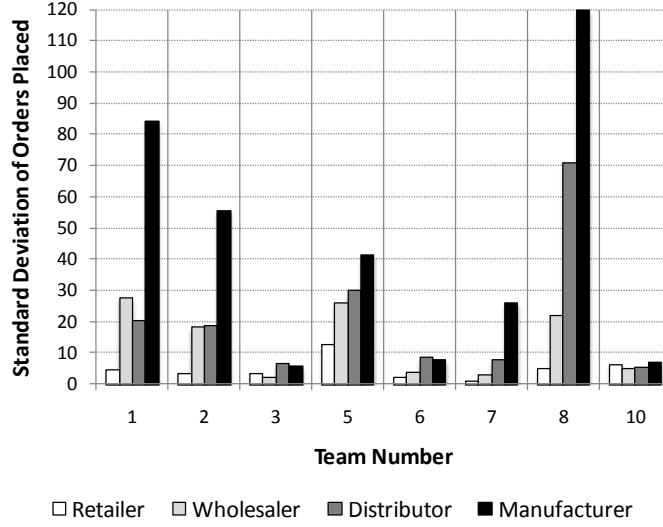


Figure 3: Standard Deviation of Orders in Experiment 1 (Baseline).

Following Sterman (1989a) we estimate the following decision rule for orders $O_{i,t}$ placed in week t by the person in role i :

$$O_{i,t} = \text{Max}\{0, CO_{i,t} + \alpha_i(S'_i - S_{i,t} - \beta_i SL_{i,t}) + \varepsilon_{i,t}\} \quad (1)$$

where CO is expected Customer Orders (orders expected from the participant's customer next period), S' is desired inventory, S is actual on-hand inventory (net of any backlog), and SL is the supply line of unfilled orders (on-order inventory). Orders are modeled as replacement of (expected) incoming orders modified by an adjustment to bring inventory in line with the target.

The parameter α is the fraction of the inventory shortfall or surplus ordered each week. The parameter β is the fraction of the supply line the participant considers. The optimal value of β is 1, since participants should include on-order as well as on-hand inventory when assessing their net inventory position. Given $\beta = 1$, the optimal value of α is also 1: since there are no adjustment costs, participants should order the entire inventory shortfall each period. The optimal expectation for customer orders in our experiment where customer demand is stationary (indeed, constant and known to all participants) is the mean of actual orders placed by the final customer, that is, $CO_{i,j} = 4$ cases/week for all sectors i . In this case the rule reduces to the familiar order-up-to rule. The parameter S' represents the sum of the desired on-hand and desired on-order inventory. Since final customer demand is constant and known, desired on-hand inventory is zero, and desired on-order inventory is the inventory level required to ensure deliveries of 4 cases/week given the order fulfillment lead time of 4 weeks (3 for the factory), yielding $S' = 16$ units (12 for the factory).

Equation (1) can also be interpreted as a behavioral decision rule based on the anchoring and adjustment heuristic (Tversky and Kahneman 1974, Kahneman *et al.* 1982), in which expected customer

orders $CO_{i,j}$ represents the anchor (order what you expect your customer to order from you), and inventory imbalances motivate adjustments above or below the forecast. To capture the possibility that participants do not use the optimal forecast of incoming orders of 4 cases/week, but rather respond to the actual orders they receive, we model expected customer orders as formed by exponential smoothing of actual incoming orders, IO , with adjustment parameter θ_i , as in Sterman (1989a),

$$CO_{i,t} = \theta_i IO_{i,t-1} + (1 - \theta_i) CO_{i,t-1} \quad (2)$$

We use nonlinear least squares to find the maximum likelihood estimates of the parameters θ_i , α_i , β_i , and S'_i , subject to the constraints $0 \leq \theta_i, \alpha_i, \beta_i \leq 1$ and $S'_i \geq 0$. We estimated the confidence intervals around each estimate with the parametric bootstrap method widely used in time series models of this type (Efron and Tibshirani 1986, Li and Maddala 1996, Fair 2003), assuming an iid Gaussian error ε with variance given by the variance of the observed residuals for each individual, σ_i^2 . An ensemble of 500 bootstrap simulations were computed for each participant, and the 95% confidence intervals for each of the four estimated parameters computed using the percentile method. Table 2 summarizes the results.

	θ	α	β	S'	R^2	RMSE
Experiment 1 (Baseline)						
Median Estimate	0.09	0.47	0.08	4.42	0.54	5.47
Median Width of 95% Confidence Interval	0.55	0.30	0.29	12.66		
N	24	32	30	30		
Experiment 2 (Common Knowledge)						
Median Estimate	0.08	0.28	0.10	8.26	0.57	2.83
Median Width of 95% Confidence Interval	0.31	0.22	0.37	10.90		
N	30	40	39	39		
Experiment 3 (Eliminating Coordination Risk)						
Median Estimate	0.00	0.24	0.14	3.48	0.29	1.62
Median Width of 95% Confidence Interval	0.91	0.30	0.50	7.76		
N	9	27	26	26		
Experiment 4 (Coordination Stock)						
Median Estimate	0.08	0.26	0.24	4.97	0.54	2.53
Median Width of 95% Confidence Interval	0.40	0.23	0.42	8.01		
N	21	28	27	27		
Median Estimates, Sterman (1989a)						
	0.25	0.28	0.30	15	0.76	2.60
N	44	44	40	40		
Median Estimates, Croson and Donohue (2002)*						
	N/A	0.22	0.14	N/A	0.71	
N		44	44		44	

Table 2: Estimated parameters for the ordering decision rule (eq. 1).

Note that the number N of estimates can differ for each parameter. N can be smaller than the number of participants in each condition (40) for four reasons. First, outliers are excluded. Second, θ cannot be identified if incoming orders are constant, which is always true for retailers and also true for upstream players if the downstream player always orders 4 cases/week (an occurrence that is more common in exp. 3 where three of the four players are automated and play optimally at all times). Third, none of the parameters can be identified when the human plays optimally by ordering 4 cases/week at all times (also more prevalent in exp. 3). Finally, when the optimal estimate of $\alpha = 0$, β and S' are undefined.

* Croson and Donohue (2002) estimated a slightly different decision rule: $O_{i,t} = \text{Max}[0, a_i S_{i,t} + b_i SL_{i,t} + c_i IO_{i,t} + d_i + \varepsilon_{i,t}]$ where IO is incoming orders and a , b , c , d are estimated. The inventory and supply line weights α and β can be expressed in terms of a and b by interpreting cIO as the participant's demand forecast and noting that the inventory correction term $d + aS + bSL = a[(d/a) - (S + (b/a)SL)]$, implying $\alpha = -a$ and $\beta = b/a$.

As expected, the uncertainty bounds around θ are large. For many participants, the variance in incoming orders is small, and sometimes zero, e.g., when their customer always orders the optimal quantity of 4 cases per week. In such cases, the smoothing time constant θ cannot be identified. More important are the estimates of α and β . In experiment 1, the median fraction of the inventory shortfall corrected each period is 0.47, and the median fraction of the supply line participants consider is 0.08, both far below the optimal values of 1. The estimates of α and β are generally tight: 97% of the estimated values of α are significantly greater than 0, implying that most participants respond to inventory imbalances, as expected. However, 72% are significantly less than the optimal value of 1, implying most participants do not use the optimal order-up-to rule but instead order only a fraction of their inventory shortfall each period. For β , 80% of the estimates are significantly less than the optimal value of one. Further, 60% are not significantly different from zero, and the best estimate of β is zero for 30% of participants, indicating that participants tend to underweight or ignore on-order inventory.

4.1.2. Discussion

Contrary to predictions based on traditional assumptions of rationality and common knowledge, the bullwhip effect remains even when demand is constant and known to all participants. Constant demand and common knowledge of it did not eliminate supply line underweighting. The median estimate of β , the fraction of on-order inventory accounted for by the participants, at 0.08, is far below the optimal estimate of 1, and estimates of desired total inventory, S' , are far too small to account for the replenishment lead-time. The median estimate of S' is 4.42 cases, while the optimal value is 16 (12 for the factory). As found in Sterman (1989a) and Croson and Donohue (2006), participants substantially underweight the supply line and underestimate the replenishment cycle time. Similar to prior studies, there is heterogeneity across participants in the estimated parameters. For example, the estimates of α and β span the full range of possible values, $[0, 1]$, and the standard deviation of the set of estimated values is large (0.35 and 0.34, respectively).

Prior work (Thomsen, Mosekilde and Sterman 1992, Mosekilde, Larsen and Sterman 1991, Sterman 1988) mapped the relationships between the parameters of the decision rule and supply chain performance. Although the parameter space is highly nonlinear, supply line underweighting (low β) and aggressive ordering (high α) can interact to increase the bullwhip, raise costs, and can produce periodic, quasiperiodic, and chaotic behavior. These studies show that, conditional on supply line underweighting, a less aggressive response to inventory discrepancies (lower α) increases stability and lowers costs by reducing the gain of the main negative feedback regulating inventory. In section 5.3 we examine whether participants in each of the four experiments, despite significantly underweighting the supply line, placed

orders in a way that approximates the optimal response to inventory shortfalls conditioned on the degree of supply line underweighting they exhibited.

Post-play questionnaire responses suggest that coordination risk may be contributing to order oscillation and amplification. A number of subjects asserted that while *they* realized that ordering 4 units every period would minimize costs, they didn't believe their *teammates* understood or would follow the optimal policy. Some explicitly note that this uncertainty caused them to seek additional inventory as protection against their teammates' unpredictable ordering behavior:

Since I was the retailer it was easy to decide what my inventory should be. I tried to keep the inventory at 4, just in case something happened along the line and I wouldn't get my supply.
[Retailer]

"I wanted to have a little extra inventory than what was demanded" [Wholesaler]

"I tried to anticipate incoming orders. [Ordering] 10... would allow an ample buffer in case a large order came in." [Distributor]

"I tried to anticipate what the rest of my teammates were going to order. This was difficult however since my teammates' orders would jump from 6 to 200 and back to 4. I tried to get inventory to help but then they stopped making order (sic) and I was screwed." [Factory]

Uncertainty about the behavior of others motivated these participants to hold additional inventory as a buffer against the risk of errors by their teammates. However, the only way a participant can increase on-hand inventory from the initial level of zero is to temporarily increase orders, thus forcing the supplier into an unanticipated backlog situation and, to the extent upstream players underweight the supply line, triggering the oscillations we observe. The questionnaire responses and the persistence of the bullwhip effect with constant, publicly known demand are consistent with the conjecture that coordination risk contributes to an increase in order oscillations and amplification.

4.2. Experiments 2 and 3: Reducing Coordination Risk

Figure 4 reports the results of experiment 2. Compared to experiment 1 (Figure 3 – which is on a larger scale), publicly providing the optimal policy appears to substantially reduce order variability. We use a nonparametric two-tailed Mann-Whitney test (also referred to as a Wilcoxon test, see Seigel 1965, p. 116) to compare the standard deviations of orders between experiments 1 and 2. The median standard deviation of orders falls significantly from the baseline value of 15.7 to 5.4 (Mann-Whitney $p = 0.012$),⁵ supporting Hypothesis 2. However, the results also show the bullwhip effect is resilient. Orders still oscillate (although oscillation is much reduced) and there remains significant amplification in order variability up the supply chain ($\sigma_i > \sigma_{i-1}$ for 83% of the cases, $p = 0.0001$).

⁵ Including outliers also yields significantly higher standard deviation of orders placed in the baseline (median 20.9) compared with those placed here (median 5.4), Mann-Whitney $p=0.007$.

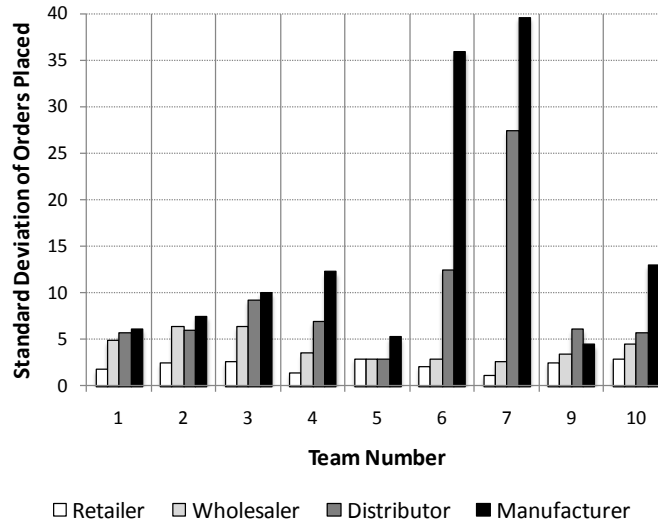


Figure 4: Order Standard Deviations in Experiment 2 (Creating Common Knowledge)

Interestingly, the tendency to underweight the supply line is not eliminated despite knowledge of the optimal policy. The median fraction of the supply line participants consider in experiment 2 is 0.10, substantially below the optimal value of 1. Further, 82% of the estimated values of β are significantly less than one, 51% are not significantly different from zero, and zero is the best estimate of β for 21% of the participants. Contrary to Hypothesis 3, providing participants with the optimal policy does not reduce supply line underweighting.

Figure 5 reports the results of experiment 3. Properly assessing the impact of automating the orders of other players (experiment 3) is complicated by the difference in the environment faced by the humans in experiment 3. Humans playing the upstream positions (W, D, or F) face no variability in incoming orders, since all downstream players are programmed to order the optimal 4 cases per week.⁶ Thus human players in these roles do not experience any unintended inventory changes, unlike their counterparts in the previous two experiments, who typically experience large variation in incoming orders. In experiments 1 and 2, retailers are the only ones guaranteed to face constant orders. The impact of automation is therefore best examined by comparing the behavior of humans in experiment 3 against the behavior of retailers in experiments 1 and 2. Doing so provides a conservative measure of the impact of the treatment, since we would expect more significant differences in the behavior of upstream players across treatments.

⁶ While the automated agents are guaranteed to play the optimal strategy, their orders may differ from 4 cases per week if the human player does not order optimally. If the human deviates from equilibrium (for any reason), upstream agents experience unintended changes in their inventory, requiring them to alter their own orders to correct the imbalance (using the optimal order-up-to rule to do so).

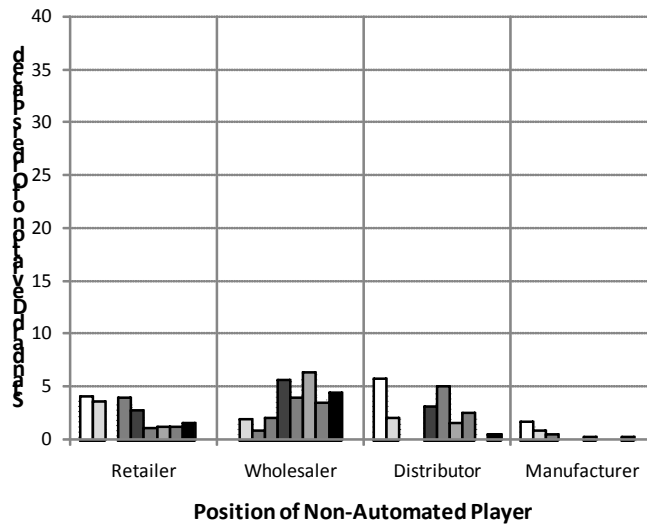


Figure 5: Order Standard Deviations in Experiment 3 (Eliminating Coordination Risk)

Note: Here participants all make up separate teams (with automated team members) and so we group the data by role rather than team.

The median standard deviation of the human players' orders in experiment 3 is 1.7, compared to 3.9 for retailers in the baseline, a significant difference (Mann-Whitney $p = 0.012$), supporting Hypothesis 4, part (a).⁷ While the median standard deviation in orders of 1.7 is lower than the value of 2.5 for retailers in experiment 2 (common knowledge), the difference is not statistically significant (Mann-Whitney $p = 0.43$).⁸ Thus the data do not consistently support part (b) of Hypothesis 4.

Nine of the forty participants in experiment 3 ordered the optimal quantity of 4 cases/week at all times. The parameters of the decision rule cannot be estimated for these subjects. For the remaining 31 participants, the median fraction of the supply line participants consider is 0.14, significantly less than the optimal value of 1. Supply line underweighting remains robust among these participants despite knowledge of the optimal ordering rule and the elimination of coordination risk: 96% of the estimates of β are significantly less than the optimal value of one. Hypothesis 5 is, at best, only weakly supported.

Overall, results from experiments 2 and 3 are consistent with the conjecture that coordination risk is one behavioral cause of the bullwhip effect. When uncertainty in the ordering behavior of one's supply chain partners is reduced (through the introduction of common knowledge of the optimal policy; experiment 2) or eliminated entirely (through the additional introduction of trust by guaranteeing others will use the optimal policy; experiment 3), order oscillations are significantly lower. Nevertheless, order

⁷ Including outliers, median standard deviations by baseline retailers is 3.9 while median standard deviations among all human players here is 1.9, Mann-Whitney $p=0.021$.

⁸ Including outliers, median standard deviations by retailers with common knowledge is 2.5, while median standard deviations here among all human players is 1.9, Mann-Whitney $p=0.27$.

oscillation and amplification are not completely eliminated, and supply-line underweighting persists.

4.3. Experiment 4: Adding Coordination Stock

The previous experiments reduce the magnitude of the bullwhip effect by reducing or eliminating coordination risk. However, in the field it is often not practical or possible to create common knowledge of the optimal policy, or to guarantee that one's supply chain partners will follow a particular ordering rule. In such cases, an intervention is needed to mitigate the impact of coordination risk, despite not reducing the risk itself. Figure 6 displays the average standard deviations across all teams for experiment 4, where excess inventory was introduced as a possible intervention. Order oscillation clearly remains, although the addition of coordination stock substantially and significantly dampens variation compared to the baseline setting (comparing Figures 3 and 6). The sign test suggests that order amplification still exists, but it is only marginally significant ($\sigma_i > \sigma_{i-1}$ in 62% of the cases, differing from the chance rate of 50% at $p = 0.097$).⁹ The median standard deviation of orders placed is 3.8 compared with 15.7 in the baseline, a significant difference (Mann-Whitney $p = 0.001$). Our results thus support hypothesis 6.

Managers and researchers recognize that inventory may be held for different reasons. Inventory held because of the need to batch orders to amortize fixed ordering costs is commonly referred to as *cycle stock*. Inventory used to buffer the effects of process uncertainty (e.g., to prevent starvation between work cells) is known as *buffer inventory*. *Safety stock* usually refers to inventory used to buffer against the impact of demand or supply uncertainty. While coordination risk is a contributing factor to demand and supply uncertainty, such uncertainty in most inventory models is considered to arise from exogenous sources. In contrast, the uncertainty subjects experience here does not arise from the external environment—it is self-inflicted. We introduce the term *coordination stock* to represent a type of safety stock used to buffer against *strategic* uncertainty in orders or delivery due to coordination risk. Coordination stock buffers against decision errors, but does not reduce those errors directly.

It is important to note that, all else equal, adding excess inventory decreases the probability of entering the backlog regime where the system is intrinsically less stable. So, even if the source of perceived risk is something other than (or in addition to) coordination risk (e.g., trusting the demand information or other aspects of the experimental setting), then adding excess inventory can still be beneficial in reducing the bullwhip effect. Also, while excess inventory reduces order oscillation, one could argue that the improvement comes at the expense of increased costs as participants hold larger inventories. Whether adding coordination stock is beneficial from a cost perspective depends on the underlying cost parameters.

⁹Teams 6, 7 and 10 were removed as outliers. Including them yields $\sigma_i > \sigma_{i-1}$ in 63% of the cases, and the p-value falls to 0.051.

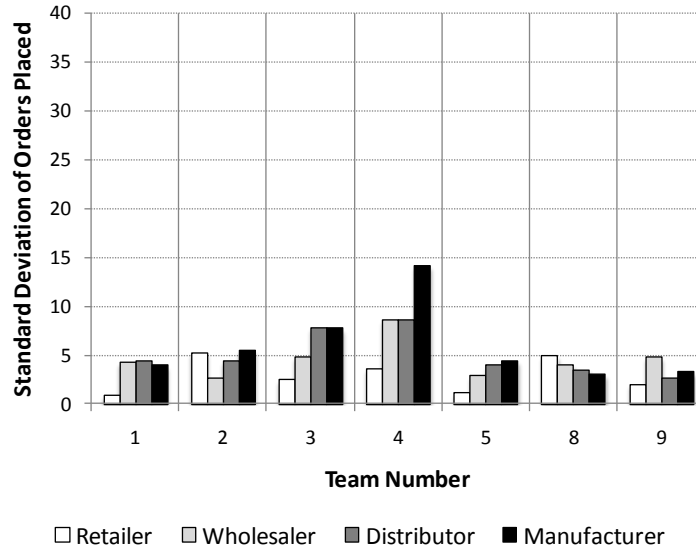


Figure 6: Standard Deviation of Orders in Experiment 4 (Adding Coordination Stock)

Specifically, one must trade off the cost of holding a larger buffer inventory (coordination stock) against the cost of increased inventory spikes and backlogs without the buffer. For our experimental setting, where the unit backlog cost is twice the unit holding cost, the median total supply chain cost with coordination stock is \$1,890, significantly lower than the median total cost of \$9,151 without coordination stock (Mann-Whitney $p = 0.0065$).¹⁰ Here, at least, the cost of holding coordination stock is more than offset by the lower costs resulting from greater stability. Our setting (where unit backlog cost is at least twice that of unit holding cost) is representative of industries where stockouts not only lead to lost sales but also erode a firm’s reputation as a reliable supplier, potentially leading to loss of market share, lower prices, and other costs. However, if unit holding costs exceed backlog cost, carrying additional inventory may be prohibitively expensive even if it reduces supply chain instability. Such a situation may exist for perishable and customized goods, products with high embodied value-added, product variants ordered infrequently, and other settings where firms tend toward make-to-order strategies.

To test the impact of excess inventory on the tendency to underweight the supply line, we again estimate the decision rule in equations (1) and (2). The median fraction of the inventory shortfall corrected each period is 0.26, and the median fraction of the supply line participants consider is 0.24, both substantially below the optimal values of 1. Furthermore, 63% of the estimates of β are significantly less than the optimal value of one, 41% are not significantly different from zero and the best estimate is zero

¹⁰ With the outlier teams the median standard deviation of orders in the baseline is 20.9 and here is 4.4, which are significantly different (Mann-Whitney $p = 0.032$). The median total supply chain cost with outliers is 13,105 in the baseline and 2,840 here, and this difference is weakly significant (Mann-Whitney $p = 0.09$).

for 19% of participants. Supply line underweighting persists. However the *consequences* of underweighting are lessened considerably by the addition of coordination stock. One reason is simply that the extra buffer means the supply chain operates farther from the unstable backlog regime.

Comparing the parameters of the decision rule between experiments 1 and 4 (Table 2) suggests another source for the benefit of coordination stock. The estimated values of α are lower in experiment 4, where participants begin with 12 units of on-hand inventory, than in experiment 1, where they begin with the optimal level of zero (a median of 0.26 vs. 0.47, respectively), although the difference is only marginally statistically significant ($p = 0.066$). Why might participants be less aggressive when given coordination stock? One possibility of course, is that, by chance, the participants in experiment 4 were calmer by disposition than those in experiment 1. Alternatively, the participants may have reacted to the situation. With the standard cost structure we use, backlogs are more costly than on-hand inventory. When initial inventory is zero, even a small increase in incoming orders creates a costly backlog, perhaps triggering an aggressive effort to return to nonnegative on-hand inventory. However, with initial inventory of 12, the same small jump in incoming orders will leave the participant with positive on-hand inventory, avoiding the costly backlog regime. Participants may therefore not feel as much urgency to order in the presence of coordination stock.

As an analogy, consider driving. You seek to maintain a certain distance between your car and the car ahead, analogous to the desired on-hand inventory level in the supply chain. Although you know that the speed limit is (say) 55 miles per hour, and although this is common knowledge, you cannot trust the driver ahead of you to maintain a steady speed of 55. You must often adjust your speed to maintain a safe distance between you and the car ahead, accelerating when the distance grows, and easing off the gas or applying the brake when the distance shrinks below the safe level. As in supply chains, the velocity of your car responds to your control actions with a lag. And, as in the beer game, costs are asymmetric: the cost of a larger gap than you desire between you and the car ahead is small, while the cost of “backlog” (rear-ending the car in front of you) is very high. Drivers who tailgate must brake aggressively when the distance to the car ahead shrinks (that is, use a large α), while drivers with a large buffer between themselves and the car ahead can simply ease off the gas (use a smaller α). That is, the gain of the negative feedback through which drivers control the distance to the car ahead is a nonlinear function of the distance. We hypothesize that, at least with the cost structure used here, participants who find themselves close to the costly backlog regime, as in experiment 1, are likely to order more aggressively than those who are given an initial buffer, as in experiment 4. The contrast between the estimated values of α in these experiments is consistent with the “tailgating” theory, but further experiments would be needed to test for the hypothesized nonlinear dependence of α on net inventory and the cost structure.

5. General Discussion and Psychological Mechanisms

The contributions of this paper are both conceptual and empirical. First, we identified coordination risk as a possible new behavioral cause of the bullwhip effect in supply chains. Second, we presented experimental evidence demonstrating that reducing factors contributing to coordination risk moderates the bullwhip effect. Third, we introduced a successful managerial intervention, showing that carrying coordination stock to buffer the system against strategic uncertainty improves performance.

However, the more basic question regarding the psychological processes through which these effects operate remains open. While identifying the psychological foundations of the phenomenon was not our objective, we can shed some light on possible mechanisms using the data we have collected.

5.1 Does the Presence of Coordination Risk Create Supply Line Underweighting?

The first possibility is that reducing coordination risk reduces supply line underweighting. Studies show that people have great difficulty managing complex dynamic systems, typically failing to account for feedback processes, underweighting time delays, and misunderstanding stocks and flows (Sterman 1994, Booth Sweeney and Sterman 2000). Underweighting the supply line of unfilled orders in stock management tasks such as the beer game is particularly common.

Sterman (1989a) discusses explanations for supply line underweighting grounded in the theory of bounded rationality and behavioral decision making (Simon 1979, Tversky and Kahneman 1974, Kahneman, Slovic and Tversky 1982). Research shows that people tend to ignore or underestimate time delays in dynamical systems (Sterman 1989a, Dörner 1996). If people's mental models do not include the delay between placing and receiving orders, they may simply be unaware that they should account for on-order inventory in order decisions, even if, as in the present experiment, that information is available. If so, providing people with the optimal policy, which explicitly instructs them to consider all on-order inventories, should eliminate, or at least reduce, supply line underweighting. Alternatively, people may know that they should take the supply line into account, but find it difficult to do so due to cognitive limitations such as limited attention, working memory, and mental computation capability (Simon 1997; Plous 1993). One must not only recognize the importance of the supply line, but direct attentional resources to monitoring it and incorporating it in the decision heuristic. If people understand the importance of the supply line, but face cognitive limits that prevent them from accounting for it, then reducing the cognitive load of the task should free up working memory and attentional resources for that purpose, reducing supply line underweighting.

The cognitive resources available to attend to the supply line should be the smallest in experiments 1 and 4, where the optimal ordering policy is not common knowledge. Cognitive resources should increase when players are given common knowledge of the optimal policy (experiment 2), and increase further

when players are guaranteed that all others in the supply chain will in fact play the optimal policy (experiment 3). Consequently, if cognitive limitations contribute to supply line underweighting, the estimated values of β should be lowest in experiments 1 and 4, higher in experiment 2, and highest in experiment 3.

To test these hypotheses we compare the estimated values of β across experiments. In experiment 1 the median value of β is extremely low (0.08), indicating that many subjects essentially ignored the supply line altogether. The median value of β in experiment 2 is nearly identical (0.10) and there is no statistically significant difference compared to experiment 1 (Mann-Whitney $p = 0.63$). Providing people with the optimal policy does not reduce supply line underweighting. Adding trust by automating all players but one (experiment 3) frees up additional cognitive resources because players know and need not speculate on how their partners will behave. However, although the median value of β rises to 0.14 the difference compared to experiment 1 is not statistically significant (Mann-Whitney $p = 0.27$), and it remains far below the optimal value of 1. The hypothesis that eliminating coordination risk by requiring others to use the optimal policy will reduce supply line underweighting is not supported. The highest estimates of β are found in the coordination stock condition (experiment 4), in which participants are not informed of the optimal policy: the median estimate of β in experiment 4 is 0.24, though these values are also not significantly different than those in experiment 1 (Mann-Whitney $p = 0.14$). Overall, there is no evidence that providing common knowledge of the optimal policy or even guaranteeing that others will use it reduces supply line underweighting.

We are left to conclude that (1) the fraction of the supply line participants take into account remains far below the optimal level in all experimental conditions, and (2) the performance improvements induced by reducing coordination risk are at best only weakly related to the degree of supply line underweighting. These results are consistent with prior work. For example, Sterman (1989b) and Diehl and Sterman (1995) show underweighting is robust in one-person stock management tasks where by definition there are no strategic interactions among players and hence no issues of coordination risk. It is possible that the complexity of those tasks still overwhelmed people's cognitive capacities, but subsequent experiments with far simpler systems show similar errors (e.g., Booth Sweeney and Sterman 2000, Cronin, Gonzalez and Sterman 2009). Prior research suggests that people suffer from deeper, more persistent difficulties in understanding and controlling dynamical systems, including difficulties understanding feedback processes, time delays, accumulations, and nonlinearities (Sterman 1994, 2011). All of these elements of dynamic complexity are present in the beer game and play important roles in its dynamics.

5.2 Does the Presence of Coordination Risk Cause Spontaneous Deviations?

The second possibility is that reducing coordination risk would reduce the tendency of individuals to

deviate from optimal ordering in the first place. Conceptually, one can distinguish between (1) the stability of the system's response to perturbations, caused by underweighting the supply line, and (2) the sources of perturbations that trigger the response.

In previous experimental studies using the beer game, exogenous perturbations such as unanticipated changes in customer demand or random variations in production and deliveries served as triggers that knocked the system out of equilibrium, allowing any biases or errors in decision making such as supply line underweighting to be observed. Here there are no such exogenous perturbations. Instead, agents who suspect their customers or suppliers will make poor decisions may endogenously choose to deviate from the equilibrium strategy to build a buffer stock against the coordination risk of non-optimal behavior by others, thus causing variability in the orders they impose on their suppliers. Note that such behavior is fundamentally different from supply line underweighting, which causes fluctuations and amplification *in response* to perturbations. We define a *spontaneous deviation* as a player ordering a quantity different from the optimal value of 4 cases *before* experiencing any change in incoming orders, deliveries, or inventory. If this mechanism drives the results, then removing the sources of coordination risk should dampen bullwhip behavior by reducing the likelihood and magnitude of spontaneous deviations.

In experiment 1, 72% of participants spontaneously deviated from the equilibrium order of 4, and 74% of these did so in the first period. Of those who spontaneously deviated, 74% ordered *more* than four cases; on average, those deviating spontaneously ordered an extra 2.6 cases above the equilibrium throughput of 4. However, the likelihood of spontaneous deviation in Experiment 2 (when coordination risk is reduced) is almost identical, 75%. Like experiment 1, the vast majority of spontaneous deviations occur in the first period (83%) and 93% of those doing so order more than four cases. The average excess order, 6.3 cases, is even larger than in experiment 1. The likelihood and magnitude of spontaneous deviations are not significantly lower in experiment 2 than in experiment 1. Spontaneous deviations are almost identical in experiment 3 (69% of participants spontaneously deviate and the average deviation is 3.9 cases), which is not significantly different than the frequency or size of spontaneous deviations in experiment 1. Thus, the likelihood and extent of spontaneous deviation cannot explain the improved performance.

5.3 Does Coordination Risk Cause Overreaction to Inventory Shortfalls?

Coordination risk might also operate by affecting people's responses to inventory shortfalls. If decision makers believe that their counterparts will make rational decisions (because those counterparts are either informed of, or constrained to follow, the optimal ordering rule), then an inventory shortfall is not a cause for alarm. Thus individuals in treatments with reduced coordination risk may be less likely to over-react to inventory shortfalls by placing excess orders.

To test this possibility, we compared the estimated values of α across the experimental treatments. The median estimate of α in experiment 1 is 0.47, while in experiment 2 (with coordination risk decreased), the median value is 0.28. While the difference is not statistically significant (Mann Whitney $p = 0.12$), the magnitude of the drop is large and has a substantial impact on costs, amplification, and oscillation.

To illustrate, simulating the system with the decision rule in equations (1)-(2) using the median estimated parameters for experiment 1 in all four positions yields a large cycle with amplitude growing over the first 48 weeks.¹¹ The first peak in orders for each link in the chain is 6, 8, 13, and 16 cases for the R, W, D and F (respectively); total costs are \$1650. Reducing α to 0.28, the median value in experiment 2—with all other parameters unchanged—cuts initial peak orders to 5, 7, 9, and 10, respectively. Peak factory orders drop 38% and total costs fall to \$1332, a drop of 19%. Thus, although the differences in the estimated values of α between experiments 1 and 2 are only directional, they are consistent with our observed results that common knowledge decreases order oscillation and amplification.

We find a similar result comparing experiments 1 and 3. The median estimate of α in experiment 3 is 0.24; the difference between this and experiment 1 is marginally significant (Mann-Whitney $p = 0.067$). Simulating the system with the median parameters from experiment 3 for one player and the optimal policy for the roles that were automated essentially eliminates amplification and oscillation, as observed in the data. In experiment 4 the median estimate of α is 0.26; again, marginally different from experiment 1 (Mann-Whitney $p = 0.09$).

Although the evidence is of marginal statistical significance, it appears that common knowledge of the optimal policy, a guarantee that others will use it, or coordination stock to buffer players against the risk of non-optimal behavior by others reduces how aggressively participants respond to a given inventory discrepancy. Since players in all conditions significantly underweight the supply line, the reduction in the gain of the inventory control feedback causes oscillation and amplification to drop. In terms of the tailgating analogy in section 5.1, a driver who fears that the cars ahead may slam on their brakes for no apparent reason must be prepared to react aggressively, while drivers who can trust that the cars ahead will not behave erratically can react less aggressively.

The persistence of the bullwhip effect, supply line underweighting, and spontaneous deviations despite the experimental treatments is consistent with recent research in judgment and decision making (e.g., Kahneman 2002, Gigerenzer *et al.* 1999) and neuroscience (e.g., McClure *et al.* 2004, Camerer *et al.* 2005) stressing the role of unconscious, intuitive processes. Imaging studies show that decision tasks involving long time frames or purely economic considerations tend to activate brain regions associated

¹¹ The system is perturbed from equilibrium because the estimated target inventory level, S^* , is not the equilibrium value. The estimated values of players' desired inventory levels induce the spontaneous deviations that knock the system out of equilibrium.

with deliberation, such as the frontal cortex, while tasks involving immediate needs and emotionally-laden contexts tend to activate the limbic system. These intuitive processes are often adaptive in settings similar to those in which they evolved, but can lead to error and bias in novel contexts. When humans evolved, complex artificial systems like supply chains did not exist. Resource scarcity placed a high premium on immediate, certain rewards compared to delayed, uncertain rewards. As shown in studies of hyperbolic discounting (e.g., McClure *et al.* 2004), people strongly prefer the “bird in the hand” and discount future rewards. In the beer game context, such preferences would lead to underweighting of the supply line of on-order inventory compared to on-hand inventory. Similarly, when humans evolved, social interactions frequently involved coordination risk: there were no automated agents guaranteed to use optimal decision rules. The supply chain context may trigger unconscious fears that others either don’t know or won’t use the optimal rule, even when people are informed of it or know that the other agents are programmed to use it. If so, the result is likely to be spontaneous deviations as players order more than the equilibrium throughput in an attempt to build coordination stock. We speculate that the robustness of supply line underweighting and spontaneous deviations in the experiment may arise in neural structures distinct from those responsible for the deliberative processes underlying rational choice. Further work investigating the psychological and neural processes at work is needed to test these hypotheses.

5.4 Does Supply Line Underweighting Affect the Aggressiveness of Inventory Corrections?

As discussed above, a participant who fully accounts for the supply line of on-order inventory ($\beta = 1$) should order the full discrepancy between desired and actual total inventory each period ($\alpha = 1$). Because there are no adjustment costs in the experiment, setting $\alpha = 1$ when $\beta = 1$ corrects any inventory gap as quickly as possible without unintended inventory overshoot. However, if participants underweight the supply line ($\beta < 1$), then large values of α cause severe instability, while smaller values reduce amplification, oscillation and costs (Thomsen *et al.* 1992, Mosekilde 1996). In control-theoretic terms, α is the gain of the negative feedback regulating inventory; the lower the value of β , the lower the gain of that feedback (the smaller the α) must be to prevent costly oscillations. Although the experimental results clearly show participants did not respond to inventory discrepancies using the joint optimal values α , $\beta = 1$, the typical value of α in each experiment is much less than one (Table 2). Are the low values of α a result of participants approximating the optimal response to inventory shortfalls conditioned on the degree of supply line underweighting they exhibited?

To shed light on this question, we computed the conditionally optimal value of α given the value of β , denoted $\alpha^*(\beta)$, for the setting of experiment 4. Values of $\alpha^*(\beta)$ were determined from simulations of the full supply chain using the decision rule in equations (1) and (2), assuming that all ordering rule

parameters other than α and β were optimal. Figure 7 shows the resulting values of $\alpha^*(\beta)$. As prior research and the discussion above suggest, $\alpha^*(\beta)$ falls sharply as β decreases.¹² Figure 7 also plots the estimated values of α and β for participants in experiment 4 (excluding outliers). The estimated values do not appear to conform to $\alpha^*(\beta)$. The correlation between α and β is small ($r = .25$) and not statistically significant ($p = .21$); including outliers yields essentially no correlation ($r = .007$, $p = .96$). The data do not support the hypothesis that participants' order decisions in experiment 4 were consistent with the conditionally optimal responsiveness to inventory shortfalls given the degree of supply line underweighting they exhibited.

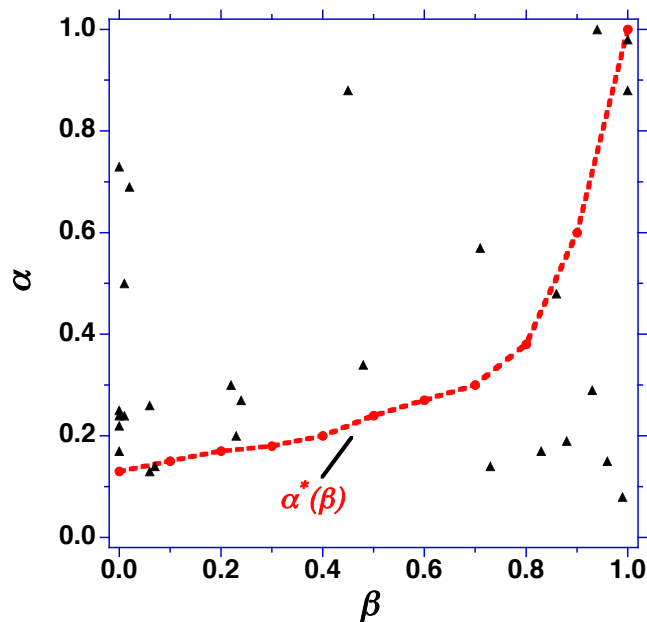


Figure 7: Estimated α and β , experiment 4, with conditionally optimal value of α as a function of β .

Conditionally optimal values of α were derived from simulations of the supply chain model under the conditions of experiment 4, using the decision rule in equations (1) and (2) and assuming optimal parameters for all links in the chain ($\theta = 0$, implying that the demand forecast is accurate at all times at 4 cases/week, and target inventory for each link is set to optimal levels of on-hand and on-order inventory). To reduce the dimensionality of the parameter space, all four simulated participants are assumed to use the same values of α and β . Estimated values shown exclude outliers.

The other experimental conditions exhibit the same pattern. Figure 8 shows the estimated (α, β) values for participants in experiments 1, 2 and 3. There is considerable scatter in all three experiments, with α spanning a wide range of values for any given value of β . The correlation between the estimated values of α and β for participants in these experiments is small ($r = .16$) and not statistically significant ($p = .12$); including outliers yields the same result ($r = .13$, $p = .16$). The evidence does not support the

¹² The exact values of the function $\alpha^*(\beta)$ plotted in Figure 7 depend on the assumptions that the other parameters in the decision rule (θ, S') are optimal and that initial inventory is 12 cases for all participants (as in experiment 4). Other assumptions and initial conditions will cause the values of $\alpha^*(\beta)$ to vary, though in all cases the conditionally optimal value of α decreases with β (Thomsen et al. 1992, Mosekilde 1996).

hypothesis that participants responded to inventory shortfalls optimally, conditional on the degree of supply line underweighting they exhibited.

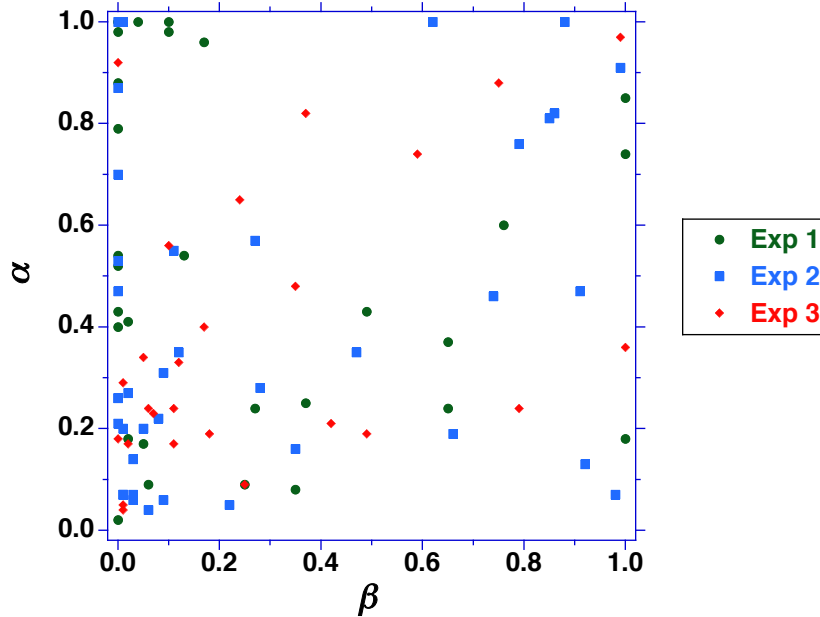


Figure 8: Estimated values of α and β in experiments 1-3 (outliers excluded).

5.5 Do Participants Engage in Hoarding Behavior?

The decision rule defined by equations (1) and (2) assumes constant desired levels of on-hand and on-order inventory. However, it is plausible that these target stock levels may vary endogenously with the state of the system. For example, when suppliers are unable to fill orders, delivery times rise and customers receive less than they desire. Customers in real supply chains sometimes respond by seeking larger safety stocks (hoarding) and by ordering more from suppliers, often through multiple channels (placing phantom orders). Hoarding and phantom orders can be rational when multiple customers compete for limited supplies, when capacity constraints constrain deliveries, and when there are supply interruptions or other stochastic shocks (Lee et al. 1997, Cachon and Lariviere 1999, Armony and Plambeck 2005). In our experimental setting, participants experience endogenous shocks due to the unanticipated ordering patterns of other supply chain members. Participants may be adjusting their target stock level in response to such shocks, effectively inducing hoarding and phantom ordering behavior. To test for the presence of hoarding and phantom ordering formally, however, requires relaxing the assumption of constant desired on-hand and on-order inventory levels in the ordering decision rule (equations (1) and (2)). For example, a participant, knowing that all members of the supply chain have

common knowledge that demand is constant, might initially attempt to maintain the optimal stock of on-hand and on-order inventory. However, if that participant received a large order, or less than they expect to receive from their supplier, they may revise their beliefs and increase the level of on-hand and on-order stocks they seek, leading to additional orders. To test this hypothesis, one of the authors of this study has estimated decision rules in which participants endogenously vary the levels of on-hand and on-order inventory they seek based on their beliefs about incoming orders, supplier lead times, and other factors indicating scarcity. The full study is available on request. In brief, the revised ordering rule is:

$$O_t = \max\left(0, D_t^e + \alpha_s (S_t^* - S_t) + \alpha_{SL} (SL_t^* - SL_t) + \varepsilon_t\right) \quad (3)$$

where target on-hand inventory, S_t^* , and target supply line of on-order inventory, SL_t^* , are functions of expected demand, D^e (defined by exponential smoothing of actual incoming orders, as in equation (2)) and expected supplier lead time, λ^e (estimated from the current supply line and the rate of deliveries received from the supplier). For example, participants may set target on-hand inventory to a certain number of weeks coverage, c , of expected incoming orders D^e :

$$S_t^* = cD_t^e \quad (4)$$

where target inventory coverage, c , may be either constant or proportional to the participant's belief about the supplier lead time. Similarly, to ensure a steady flow of deliveries from the supplier, participants must, following Little's Law, set the target supply line of on order inventory equal to λ^e weeks worth of their forecast of incoming orders,

$$SL_t^* = \lambda_t^e D_t^e . \quad (5)$$

Full exposition of the models, results and sensitivity analysis is beyond the scope of the present paper. In brief, for approximately 80% of participants the modified model does not significantly improve explanatory power. However, for the remaining 20% there is strong evidence of hoarding and/or phantom ordering: the hypothesis that target on-hand and on-order stock levels are constant is rejected and the enhanced model significantly improves R^2 and lowers the root mean square error (RMSE) between the model and data. The improvement is particularly large for those participants classified here as outliers.

Figure 9 illustrates two examples: the wholesaler and distributor from the most extreme outlier team in the data, Team 4 in Experiment 1 (shown in Figure 2). Although the retailer's orders never exceeded 36 cases/week, the wholesaler's orders reached a peak of 30,000 cases in a single week; the distributor's orders peak at 20,000 cases/week. The original decision rule with constant target stock levels cannot explain such extreme behavior. However, incorporating hoarding and phantom ordering into the model significantly improves the fit to the data, boosting R^2 from less than 0.05 to more than 0.70, and reducing the RMSE by 46% for the wholesaler and 59% for the distributor. The estimated parameters suggest

these participants reacted to the jump in incoming orders by increasing target levels for on-hand inventory and by increasing the target supply line of on-order inventory as their supplier’s lead time lengthened.

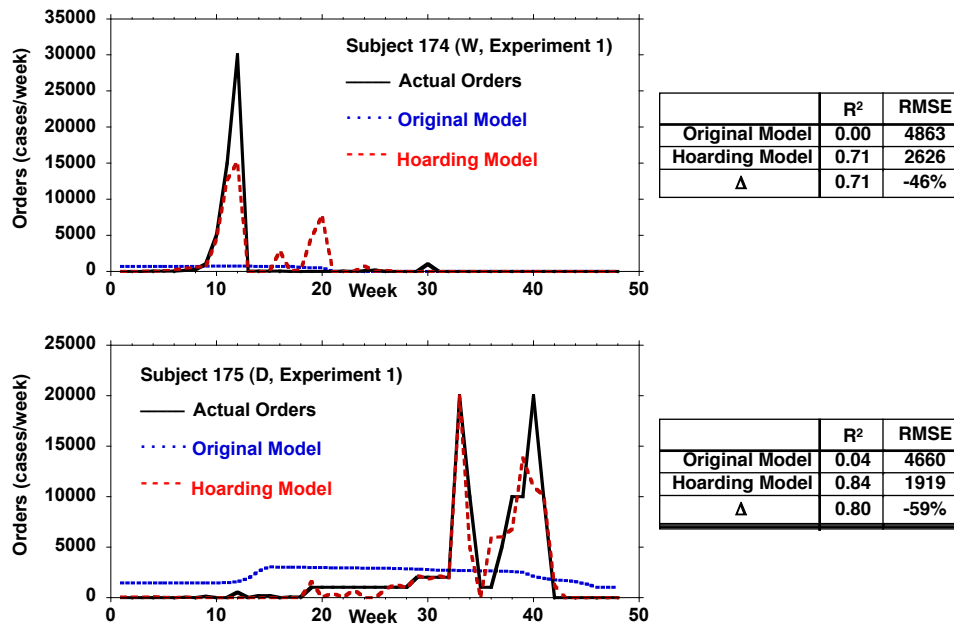


Figure 9. Fit to participant behavior of modified decision rule incorporating hoarding and phantom orders.

Why did a fifth of the participants engage in hoarding and phantom ordering when it is not rational to do so in the experiments reported here? Following the discussion in section 5.3, we conjecture that these behaviors are adaptive in situations characterized by scarcity and competition for resources, such as those in which humans evolved, and may be embedded in neural structures in the limbic system, such as those that appear to lead to hyperbolic discounting. We speculate that hoarding and phantom ordering may be triggered by the stress caused by unexpectedly large incoming orders or lags in supplier deliveries even though such behavior is not rational in the experimental setting. Further research is needed to test these hypotheses.

6 Limitations and Conclusions

Table 3 summarizes the results of our four experiments, showing which differences in order oscillation and other measures of interest are statistically significant.

Four main findings emerge from the experiments.

1. The “bullwhip effect” and supply line underweighting persist even when demand is constant and known to all (i.e., when there is common knowledge of demand). This controlled setting provides stronger evidence, as compared with previous experimental research, that the bullwhip effect is, in part, a behavioral phenomenon.

2. We identify *coordination risk* as a possible new source of uncertainty that may both trigger and further amplify bullwhip behavior. We test the impact of two factors that may reduce coordination risk, a lack of common knowledge of the optimal strategy and a lack of trust that others will follow the optimal ordering policy.

3. Performance is improved through the introduction of common knowledge and the addition of trust, although these factors do not reduce the behavioral tendency to underweight the supply line.

4. Performance is similarly improved by adding *coordination stock*, which buffers against strategic uncertainty.

Building on our notion of coordination risk, Su (2008) provides an analytical proof that such risk can lead to bullwhip behavior when the decision maker is boundedly rationality. Su's result holds when there is no information lag and orders are filled immediately (i.e. without the potential for error introduced by supply line underweighting). These results are consistent with our experimental findings and together provide strong evidence that coordination risk can inhibit supply chain performance.

Measures of Orders Placed	Experiment 1: Baseline (Retailers only)	Experiment 2: Creating Common Knowledge (Retailers only)	Experiment 3: Eliminating Coordination Risk	Experiment 4: Adding Coordination Stock	Significance tests: ** p < 0.05 * p < 0.10
Outlier teams excluded					
Median std deviation	15.7 (Retailers =3.9)	5.4 (Retailers = 2.5)	1.7	3.8	1 > 2** 1 > 3**
Average std deviation	20.6 (Retailers = 4.7)	8.3 (Retailers = 5.5)	1.9	4.5	2 = 3 1 > 4**
Median Total Supply Chain Cost	9,151 (Retailers =336)	3,396 (Retailers = 122)	47	1,890	1 > 2* 1 > 3** 2 > 3*
Average Total Supply Chain Cost	10,771 (Retailers =546)	4,562 (Retailers = 202)	147	2,157	1 > 4**
Outlier teams included					
Median std deviation	20.9 (Retailers =3.9)	5.4 (Retailers = 2.5)	1.9	4.4	1 > 2** 1 > 3** 2 = 3
Average std deviation	455.1 (Retailers =4.5)	8.3 (Retailers = 5.5)	2.3	57.8	1 > 4**
Median Total Supply Chain Cost	13,105 (Retailers =279)	3,396 (Retailers = 122)	54	2,840	1 > 2** 1 > 3** 2 > 3*
Average Total Supply Chain Cost	202,646 (Retailers =463)	4,562 (Retailers = 202)	159	14,101	1 > 4**

Note: The numbers in parenthesis for Experiments 1 and 2 are for retailers only, and are used to compare with Experiment 3.

Note: To derive these summary statistics we first calculated each individual's standard deviation of orders placed. We then calculated the median (average) standard deviation of orders placed separately for each role (R, W, D, F). Finally, we took the average of these medians (averages) to report as the summary statistic in the table and the text of the paper.

Table 3: Supply chain costs and metrics of order amplification

All research has limitations and ours is no exception. These limitations also suggest directions for future research. First, we used students as participants, as is common practice in experimental economics. It is possible that experienced supply chain professionals would perform better, although there is no evidence of any systematic differences due to the subject pool. For example, Croson and Donohue (2006) report beer game results from supply chain managers and find similar patterns of bullwhip and underweighting of the supply line. Croson (2007) provides an overview of research from experimental economics and other fields on the use of student versus professional subjects. She reports only a very few settings where professional participants perform significantly differently than student participants.

Croson and Donohue (2006), and Wu and Katok (2006) provide some initial insight into the effects of learning from experience. Taken as a whole, the evidence suggests that experience is unlikely to eliminate the bullwhip effect by itself. Managers who experience instability may conclude that others in their supply chain are unreliable, motivating them to deviate sooner or more from optimal play so as to build up coordination stock as a buffer. Further research on the role of experience, learning, and decision aids (e.g., Wu and Katok 2006) is clearly needed.

The specific choice of initial conditions and parameters suggests other extensions. Although there is no *a priori* reason to think so, it is possible that different cost parameters or initial inventory levels might lead to different outcomes.

As described above, the decision rule estimated here assumes constant levels for target on-hand and on-order inventory. These assumptions can be relaxed to test whether at least some participants endogenously increased the levels of coordination stock they sought to accumulate when faced with a surge in incoming orders or a drop in supplier delivery reliability. Preliminary work suggests that some participants in the experiment did so, and that a model incorporating behavioral responses to order volatility and supply disruptions explains the behavior of many participants classified as outliers. Further exploration of the conditions that trigger such hoarding, and empirical investigation into whether such behavioral responses to scarcity and volatility exist in real supply chains, offers a promising avenue for further research.

Finally, experiment 3 provides a starting point to investigate how people interact with artificial agents programmed to follow specific decision rules. While our agents were programmed to follow the cost-minimizing decision rule, further research could explore participants' responses to agents programmed to use other rules, including suboptimal rules or rules estimated from the behavior of human players.

Our work suggests three managerial insights. First, we demonstrate that the bullwhip effect is a behavioral phenomenon as well as an operational one, and therefore methods for reducing supply chain instability should address the behavioral as well as structural causes of the problem. Decision makers have a difficult time controlling systems that include feedbacks and delays, so training and decision aids

that help them do so may lead to substantial improvement.

Second, the notion of “optimal” behavior is contingent on people’s assumptions about the thinking and behavior of the other agents with whom they interact. If a person believes that their counterparts will behave in an unpredictable and capricious fashion, this may lead to further instability in the supply chain. To the extent that managers can be assured of their supply chain partners’ knowledge of the optimal decision rule and trust them to implement it, performance can improve further.

Our final contribution is the preliminary identification of an effective mechanism to moderate the bullwhip effect: coordination stock. The results suggest a previously unidentified purpose for excess stock. Intuitively, the optimal level of coordination stock depends on the level of coordination risk, the cost of holding excess inventory, and the cost of not dampening the bullwhip effect. Additional experiments would be needed to further define the nature of these relationships. This represents a fertile area for future research.

References

- Armony, M., E. Plambeck. 2005. The impact of duplicate orders on demand estimation and capacity investment. *Management Science*, **51** (10), 1505-1518.
- Booth Sweeney, L. and J. Sterman. 2000. Bathtub dynamics: Initial results of a systems thinking inventory. *System Dynamics Review*, **16**(4), 249-294.
- Brehmer, B. 1992. Dynamic decision making: Human control of complex systems. *Acta Psychologica*. **81** 211-241.
- Cachon, G. and C. Camerer. 1996. Loss-Avoidance and Forward Induction in Experimental Coordination Games, *The Quarterly Journal of Economics*, 111(1), 165-194.
- Cachon, G., T. Randall and G. Schmidt. 2007. In Search of the Bullwhip Effect. *Manufacturing & Service Operations Management*, **9** (4), 457-479.
- Camerer, C., G. Loewenstein, D. Prelec. 2005. Neuroeconomics: How neuroscience can inform economics. *Journal of Economic Literature* **43**, 9-64.
- Clark, A. and H. Scarf. 1960. Optimal policies for a multi-echelon inventory problem. *Management Science*. **6** 475-490.
- Chen, F. 1999. Decentralized supply chains subject to information delays. *Management Science*. **45** 1016-1090.
- Chen, F., Z. Drezner, J. Ryan, and D. Simchi-Levi. 2000. Quantifying the bullwhip effect: The impact of forecasting, lead times, and information. *Management Science*. **46**(3) 436-443.
- Chen, L. and H. Lee. 2011. Bullwhip Effect Measurement and its Implications, working paper, Fuqua School of Business, Duke University.
- Cohen, M. and M. Baganha. 1998. The stabilizing effect of inventory in supply chains. *Operations Research*, supplement to **46**(3) 572-583.
- Cronin, M., C. Gonzalez, J. Sterman (2009). Why Don’t Well-Educated Adults Understand Accumulation? A Challenge to Researchers, Educators, and Citizens. *Organizational Behavior and Human Decision Processes* 108(1): 116-130.
- Crosan, R. 2007. The use of students as participants in experimental research. *Working paper*,

- University of Texas at Dallas (prepared for Behavioral Dynamics in Operations Management Online Discussion Group).
- Croson, R. and K. Donohue. 2002. Experimental economics and supply-chain management. *Interfaces*, **32**(5) 74-82.
- Croson, R. and K. Donohue. 2003. Impact of POS data sharing on supply chain management: An experimental study. *Production and Operations Management*. **12**(1) 1-11.
- Croson, R. and K. Donohue. 2006. Behavioral causes of the bullwhip and the observed value of inventory information. *Management Science*. **52**(3) 323-336.
- Diehl, E. and J. Sterman. 1995. Effects of feedback complexity on dynamic decision making. *Organizational Behavior and Human Decision Processes*. **62**(2) 198-215.
- Dörner, D. 1980. On the difficulties people have in dealing with complexity. *Simulations and Games*. **11**(1) 87-106.
- Dörner, D. 1996. *The Logic of Failure*. Metropolitan Books/Henry Holt, New York.
- Efron, B. and R. Tibshirani. 1986. Bootstrap methods for standard errors, confidence intervals, and other measures of statistical accuracy. *Statistical Science*. **1**(1) 54-75.
- Fair, R. 2003. Bootstrapping macro econometric models, *Studies in Nonlinear Dynamics and Econometrics*. **7**(4) article 2.
- Forrester, J. 1961. *Industrial Dynamics*. MIT Press, Cambridge, MA.
- Geanakoplos, J. 1992. Common knowledge. *Journal of Economic Perspectives*. **6**(4) 53-82.
- Gigerenzer, G., Todd, P. et al.: 1999, *Simple Heuristics that Make Us Smart*. New York: Oxford University Press.
- Grubbs, F. 1969. Procedures for detecting outlying observations in samples. *Technometrics*. **11**(1) 1-21.
- Kahneman, D., Slovic, P. and Tversky, A. 1982. *Judgment under Uncertainty: Heuristics and Biases*. Cambridge UK: Cambridge University Press.
- Kampmann, C. and J. Sterman. 1998. Do markets mitigate misperceptions of feedback in dynamic tasks? Working paper, Sloan School of Management, MIT.
- Kleinmuntz, D. 1985. Cognitive heuristics and feedback in a dynamic decision environment. *Management Science*. **31** 680–702.
- Kleinmuntz, D., & Thomas, J. 1987. The value of action and inference in dynamic decision making. *Organizational Behavior and Human Decision Processes*. **39** 341–364.
- Lee, H., V. Padmanabhan, and S. Whang. 1997. Information distortion in a supply chain: the bullwhip effect, *Management Science*, **43**(4) 546-558.
- Lee, H., V. Padmanabhan, and S. Whang. 2004. Comments on Information distortion in a supply chain : the bullwhip effect, *Management Science*. **50**(12) 1887-1893.
- Lei, V., Noussair, C., and Plott, C. 2001. Nonspeculative bubbles in experimental asset markets: Lack of common knowledge of rationality vs. actual irrationality. *Econometrica*. **69**(4) 831-859.
- Li, H. and Maddala, G. 1996. Bootstrapping time series models. *Econometric Reviews*. **15**(2) 115-158.
- McClure, S., Laibson, D., Loewenstein, G., Cohen, J. 2004. Separate Neural Systems Value Immediate and Delayed Monetary Rewards. *Science*. 306 (15 October) 503-507.
- Mosekilde, E. 1996. *Topics in nonlinear dynamics: applications to physics, biology, and economic systems*. World Scientific, River Edge, NJ.

- Mosekilde, E., E. Larsen, J. Sterman (1991). Coping with complexity: Deterministic chaos in human decision making behavior. *Beyond Belief: Randomness, Prediction, and Explanation in Science*. J. L. Casti and A. Karlqvist (eds). Boston, CRC Press: 199-229.
- Nagel, R. 1995. Unraveling in guessing games: An experimental study. *American Economic Review*. **85**(5) 1313-26.
- Plous, S. 1993. *The psychology of judgment and decision making*. New York: McGraw-Hill.
- Schwieren, C. and Sutter M. 2008. Trust in Cooperation or Ability? An Experimental Study on Gender Differences. *Economics Letters*, **99**(3), 494-497.
- Seigel, S. 1965. *Nonparametric Statistics for the Behavioral Sciences*. Wiley, New York.
- Simon, H. 1979. Rational Decision Making in Business Organizations. *American Economic Review*, **69**(4), 493-513.
- Simon, H. 1997. *Administrative behavior*, 4th ed. The Free Press, NY.
- Steckel J., S. Gupta and A. Banerji. 2004. Supply chain decision making: Will shorter cycle times and shared point-of-sale information necessarily help? *Management Science*. **50**(4) 458-464.
- Sterman, J. 1988. Deterministic Chaos in Models of Human Behavior: Methodological Issues and Experimental Results. *System Dynamics Review* **4**: 148-178.
- Sterman, J. 1989a. Modeling managerial behavior: Misperceptions of feedback in a dynamic decision making experiment. *Management Science*. **35** 321-339.
- Sterman, J. 1989b. Misperceptions of feedback in dynamic decision making. *Organizational Behavior and Human Decision Processes*. **43**(3) 301-335.
- Sterman, J. 1994. Learning in and about complex systems. *System Dynamics Review*. **10**(2-3) 291-330.
- Sterman, J. 2000. *Business Dynamics: Systems Thinking and Modeling for a Complex World*. McGraw-Hill, New York.
- Sterman, J. 2011. Communicating Climate Change Risks in a Skeptical World. *Climatic Change* **108**: 811-826.
- Su, X. 2008. Bounded rationality in newsvendor models. *Manufacturing & Service Operations Management*, **10**(4) 566-589.
- Thomsen, J., E. Mosekilde, J. Sterman (1992). Hyperchaotic Phenomena in Dynamic Decision Making. *Journal of Systems Analysis and Modeling Simulation (SAMS)* **9**: 137-156.
- Tversky, A. and Kahneman, D. 1974. Judgment under uncertainty: heuristics and biases. *Science*, **185**, 1124-1131.
- Vanderschraaf, P., Sillari, G. 2005. Common knowledge, *The Stanford Encyclopedia of Philosophy* Edward N. Zalta (ed.), <http://plato.stanford.edu/archives/win2005/entries/common-knowledge/>.
- Van Huyck, J., R. Battalio and R. Beil. 1990. Tacit Coordination Games, Strategic Uncertainty, and Coordination Failure, *American Economic Review*, March 1990, 234-248.
- Van Huyck, J., R. Battalio and R. Beil. 1991. Strategic Uncertainty, Equilibrium Selection, and Coordination Failure in Average Opinion Games, *The Quarterly Journal of Economics*, **106**(3), 885-911.
- Van Huyck, J., R. Battalio and R. Beil. 1993. Asset Markets as an Equilibrium Selection Mechanism: coordination failure, game form auctions, and tacit communication. *Games and Economic Behavior*, **5**(3), 485-504.
- Wu, D. and E. Katok. 2006. Learning, Communication and the Bullwhip Effect. *Journal of Operations Management*. **24** 839-850.
- Zipkin, P. 2000. *Foundations of Inventory Management*. McGraw-Hill, New York.

Appendix: Outlier Analysis

To test for outliers, we used Grubb’s procedure, also known as the maximum normed residual test (Grubbs 1969), as recommended in the Engineering Statistics Handbook (National Institute of Standards and Technology; www.nist.gov). This procedure involves calculating the variance of orders placed over the entire game for each individual and comparing these variances to those of others in the same role/experiment category (e.g., all retailers within experiment 1). The Grubbs procedure then identifies outliers within each role/experiment category based on simple criteria; distance of an observation from the mean in terms of standard deviation. In our case, an individual was flagged as an outlier if its variance of orders was more than 2.29 standard deviations beyond the mean variance (which is consistent with a statistical significance of $p=.05$). For example, out of the 10 retailers in experiment 3, one had a variance of orders placed of 37.36. This is 2.39 standard deviations beyond the mean variance, and so this individual was identified as an outlier. For the all-human experiments (i.e. experiments 1, 2 and 4), a team was eliminated if two or more of its members did not pass the criterion. For the experiments with one human and three automated team members (i.e., experiment 3), a team was eliminated if its single subject did not pass the criterion. Table A.1 lists the individuals identified as outliers, and the number of teams eliminated as a result, for each experiment.

	Experiment 1: Baseline	Experiment 2: Creating Common Knowledge	Experiment 3: Eliminating Coordination Risk	Experiment 4: Adding Coordination Stock
Outliers by Role (team #)	Retailer (5) Wholesaler (4) Distributor (4) Factory (4) Distributor (9) Factory (9)	Retailer (8) Distributor (7)	Retailer (18) Wholesaler (27) Factory (1) Factory (5)	Distributor (6) Factory (6) Retailer (7) Wholesaler (7) Distributor (7) Factory (7) Retailer (10) Wholesaler (10)
Eliminated Teams	Teams 4, 9	None	Individuals 1, 5, 18, 27	Teams 6, 7, 10

Note: The ID number/team number of each outlier player is listed in parenthesis.

Table A.1. Breakdown of Individual Outliers and Eliminated Teams