



Management Science

Publication details, including instructions for authors and subscription information:
<http://pubsonline.informs.org>

Principal-Agent Settings with Random Shocks

Jared Rubin, Roman Sheremeta

To cite this article:

Jared Rubin, Roman Sheremeta (2016) Principal-Agent Settings with Random Shocks. Management Science 62(4):985-999.
<http://dx.doi.org/10.1287/mnsc.2015.2177>

Full terms and conditions of use: <http://pubsonline.informs.org/page/terms-and-conditions>

This article may be used only for the purposes of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval, unless otherwise noted. For more information, contact permissions@informs.org.

The Publisher does not warrant or guarantee the article's accuracy, completeness, merchantability, fitness for a particular purpose, or non-infringement. Descriptions of, or references to, products or publications, or inclusion of an advertisement in this article, neither constitutes nor implies a guarantee, endorsement, or support of claims made of that product, publication, or service.

Copyright © 2015, INFORMS

Please scroll down for article—it is on subsequent pages



INFORMS is the largest professional society in the world for professionals in the fields of operations research, management science, and analytics.

For more information on INFORMS, its publications, membership, or meetings visit <http://www.informs.org>

Principal–Agent Settings with Random Shocks

Jared Rubin

Argyros School of Business and Economics and Economic Science Institute, Chapman University,
Orange, California 92866, jrubin@chapman.edu

Roman Sheremeta

Weatherhead School of Management, Case Western Reserve University, Cleveland, Ohio 44106;
and Economic Science Institute, Chapman University, Orange, California 92866, rshereme@gmail.com

Using a gift-exchange experiment, we show that the ability of reciprocity to overcome incentive problems inherent in principal–agent settings is greatly reduced when the agent's effort is distorted by random shocks and transmitted imperfectly to the principal. Specifically, we find that gift exchange contracts without shocks encourage effort and wages well above standard predictions. However, the introduction of random shocks reduces wages and effort, regardless of whether the shocks can be observed by the principal. Moreover, the introduction of shocks significantly reduces the probability of fulfilling the contract by the agent, the payoff of the principal, and total welfare. Therefore, our findings demonstrate that random shocks place an important bound on the ability of gift exchange to overcome principal–agent problems.

Data, as supplemental material, are available at <http://dx.doi.org/10.1287/mnsc.2015.2177>.

Keywords: gift exchange; principal–agent model; contract theory; reciprocity; effort; shocks; laboratory experiment

History: Received November 6, 2013; accepted January 30, 2015, by John List, behavioral economics. Published online in *Articles in Advance* August 19, 2015.

1. Introduction

This paper addresses two related sources of inefficiency that can arise in principal–agent relationships. First, a large literature notes that if the agent's effort is signaled imperfectly to the principal and monitoring is expensive or impossible, then it may be impossible to write a first-best contract, since the observed outcomes are not perfectly correlated with the agent's actions.¹ Second, if contracts are not exogenously enforceable or verifiable, endogenous enforcement through incentive compatibility requirements generally incentivizes agents to provide suboptimal levels of effort.² These two problems are related because it is impossible to exogenously enforce a contract (through legal or other institutions) that specifies effort requirements when effort is unobservable.

The unobservable effort problem is a common one for firms, because there are many types of tasks in which effort is positively correlated with observable outcomes, but these outcomes are also a function of random shocks (such as profits, number of sales, etc.).

For example, the quantity of sales made by regional salespeople reflects both their effort and local demand fluctuations, where the latter are ostensibly random and difficult to observe. Hence, an employee can put in very little effort but perform well because of luck. Under these conditions, what is fair remuneration?³ Should the employee be punished for lack of effort or rewarded for a good performance that predominantly came from luck? On the other hand, another employee can put forth very high effort but perform poorly because of bad luck. In that case, should the employee be punished for a bad outcome or rewarded for a high effort? Despite settled theoretical predictions, there is very little empirical research investigating how luck and effort play in remuneration in settings where effort is unobservable (Charness and Kuhn 2011). This is understandable because it is difficult to measure empirically to what

¹ See, for example, Harris and Raviv (1979), Holmström (1979), Shavell (1979), Holmström and Milgrom (1991), and Baker (1992). Prendergast (1999) provides a more general overview of the contracts literature that emerged in the 1970s–1990s.

² See, for example, Grossman and Hart (1983), Milgrom and Roberts (1992), and Laffont and Martimort (2002).

³ According to the “informativeness principle” of Holmström (1979), when perfect information is not available, any observable measure of performance reveals information about the effort level chosen by the agent and should be used in the compensation contract. When effort is perfectly observable, the problem of optimal contract design is trivial: remuneration should be based on effort and not luck. This is sometimes referred to as the “accountability principle” (Konow 2000, 2003), which states that remuneration should be based on the relevant variables that an individual can influence (i.e., effort) but not those that an individual cannot influence (i.e., luck).

degree effort versus luck impacts individual performance. It is even more difficult to evaluate how employers reward effort versus luck, because remuneration is usually based on final performance, which is a function of effort, ability, and luck (Ericsson and Charness 1994).

The second problem that firms can face when contracting with employees is contract enforceability and verifiability. Even where legal institutions exist, writing a first-best, fully contingent contract is often impossible when effort is not verifiable. This problem is especially stark when random shocks affect the mapping from effort to outcome. For example, if a contract offers a wage in return for the first-best effort, the agent has incentive to provide less than the first-best effort if there is a high enough probability of getting lucky (as a result of a positive production shock). Since the principal cannot verify whether the outcome is due to effort or luck, the principal cannot enforce the contract.

Fehr et al. (1997) provide experimental evidence that the contract enforceability problem is partially mitigated by behavioral concerns for reciprocity. They build on an extensive literature that suggests that the reciprocity motivation can help to explain a host of results that are contrary to standard economic theory.⁴ One implication of this literature is that contracts based on reciprocity come closer to the first-best than standard contract theory dictates. Fehr et al. (1997) test this implication with a gift-exchange experiment, where principals offer contracts that include wages and desired effort levels. Agents who accept the contracts receive the wage and choose an effort level (where higher effort improves the principal's payoff), but they do not have to abide by the desired effort level in the contract. The money-maximizing Nash equilibrium is for the agent to provide zero effort (since it is costly and they cannot be punished) and for the principal to thus offer the lowest possible wage. In their experiment, however, agents frequently show positive reciprocity; not only do they provide more effort than the money-maximizing Nash equilibrium prediction, but their effort level is increasing in the wage offered by the principal. These results are exacerbated when principals are also allowed to exhibit reciprocity. In one treatment, Fehr et al.

(1997) introduce a third stage in which principals can pay to punish or reward the agent after observing their effort. Although the addition of this stage does not alter the money-maximizing Nash equilibrium predictions of wage or effort, they find that allowing both sides to exhibit reciprocity significantly increases effort (and thus efficiency), and that both principals and agents are better off than they are when only agents are allowed to show reciprocity. Fehr et al. (2007) provide further evidence that this type of bonus contract vastly outperforms standard incentive-based contracts despite relying on unenforceable actions.

These papers contribute significantly to our knowledge of how behavioral incentives encourage contract enforcement in the absence of explicit incentives. Yet, Fehr et al. (1997, 2007) only consider how reciprocity improves contract efficiency under perfect information. In their experiments, principals can reward or punish agents based on perfectly observed effort—there are no random shocks affecting the mapping from effort to outcome. This is an important omission, because the types of contracts they are concerned with are often difficult to enforce in the real world precisely because outcomes are affected by shocks and thus optimal effort levels are impossible to induce in an exogenously enforced contract. Indeed, it is not clear *ex ante* how the introduction of shocks interacts with the reciprocity motive. Do principals exhibit reciprocity when they are unsure that the outcome that they observe is the result of the agent's effort?

This paper addresses this problem. We conduct a gift-exchange experiment similar to that of Fehr et al. (1997), except that the principal receives an imperfect signal of the agent's effort in some treatments. Our first treatment is similar to the bonus treatment in Fehr et al. (1997). Principals and agents are randomly matched, and the principal offers a wage and asks for a desired effort. The agent then receives the wage and can choose any effort (where the cost of effort is increasing in effort chosen). The principal can then reward or punish the agent, although either is costly. There are no shocks in this treatment, so we employ it as our baseline. The second treatment is exactly the same as the first, except that we add a random (uniformly distributed) number to the agent's effort. In this treatment, there is still perfect information; the principal observes both the effort level *and* the random number when making the decision of how much to punish or reward the agent. The final treatment is exactly like the second treatment, except that principals only observe the outcome (effort + random number), not the agent's effort. Relationships in all treatments are one-shot and anonymous, so reputational concerns are absent.

⁴ There is a wealth of experimental evidence that both positive and negative reciprocity have important effects on actions, with negative reciprocity being shown as more salient. In the context of the gift-exchange experiment employed in this paper, see Charness and Haruvy (2002), Charness (2004), Charness and Dufwenberg (2006), Fehr and Schmidt (2007), and Houser et al. (2008). Rabin (1993) provides the canonical model introducing reciprocity into game theory, and Falk and Fischbacher (2006) provide a theory connecting the reciprocity motive to a host of standard experimental results. Fehr and Gächter (2000) provide a survey of the literature on reciprocity.

Consistent with previous literature on gift exchange (Fehr et al. 1997, 2007; Charness and Kuhn 2011), we find that bonus contracts without shocks encourage effort and wages well above standard predictions. However, we also find evidence that this result is partially mitigated when random shocks are present. The mere introduction of shocks reduces wages and effort, regardless of whether the shocks are observed by the principal.

What can explain our findings? Why should the introduction of a shock reduce wages and effort if the shock is perfectly observable? To address this question, we outline a model of reciprocity (in the context of our experiment) where subjects reciprocate based on either the effort or outcome of the previous game play.⁵ If wages and effort are solely encouraged by effort-based reciprocity, there should be no difference between the baseline treatment (without shocks) and the treatment where shocks are perfectly observed, since the reciprocity motive is based on the other's action, not the outcome emanating from the action. On the other hand, if wages and effort are solely encouraged by outcome-based reciprocity, there should be no difference between the treatment where shocks are perfectly observed and the treatment where shocks are not observed, since the mapping from effort to expected outcome (and reciprocity) is the same in both cases. Moreover, the effort exerted in these two treatments should be lower than in the baseline treatment, since it is more costly to "make up" for a bad shock than it is to scale back effort for a "good shock" (the cost curve is increasing and convex).

We find evidence in favor of subjects exhibiting outcome-based reciprocity. In the treatment where principals observe both the agent's effort and the shock, principals do indeed vary their adjustments based on the shock, which is outside the control of the agent. This suggests that reciprocity is in part influenced by the outcome the principal receives, even if the principal knows that part of this outcome was influenced by luck. Moreover, as the outcome-based reciprocity hypothesis suggests, wages and effort are significantly lower in treatments where shocks are perfectly observed relative to the baseline, and we observe no differences in behavior between treatments where shocks are present and observable and treatments where shocks are present and unobservable.

These results have a number of important implications. First, our results provide evidence that the

reciprocity motive in the principal-agent settings is based on the outcome (which is a function of effort and shocks) of others' actions and not simply on their intentions (i.e., effort). In this regard, we contribute to the literature studying how individual behavior is impacted by intentions and outcomes (Falk and Fischbacher 2006). We show that this result has significant welfare implications: welfare-enhancing effort is lower in the presence of shocks, even when the shocks are perfectly observable. Second, our results contribute to the large literature on gift exchange. Charness and Kuhn (2011) review the experimental evidence on gift exchange, concluding that gift exchange is a robust phenomenon in that higher wages lead to higher effort. Our study contributes to this literature by demonstrating that the existence of random shocks is an important boundary condition of gift exchange. To this end, our study adds to an important literature, highlighted by List (2007), examining how the introduction of realistic elements and institutions into gift-exchange settings impacts individual behavior (e.g., Gneezy and List 2006, Falk and Kosfeld 2006, Charness and Gneezy 2008). Gneezy and List (2006), for example, show that positive reciprocity effects detected in laboratory gift exchange experiments can wear off quickly in the field. Similarly, Falk and Kosfeld (2006) show that reciprocity declines when principals try to control agents' performance. On the other hand, Charness and Gneezy (2008) show that agents are more reciprocal when anonymity is reduced. Our results add yet one more realistic element, showing that the ability of reciprocity to overcome incentive problems inherent in principal-agent settings is greatly reduced when the agent's effort is distorted by random shocks and transmitted imperfectly to the principal (as is usually the case in the real world).

2. Experimental Design and Procedures

Our experimental design is built on a variation of a gift-exchange game. The game consists of three stages. In stage 1, the principal offers contract (w, e) to the agent; i.e., the principal offers a wage w (any integer number between 1 and 100) and the desired effort e (an integer number between 0 and 14) that the principal would like the agent to undertake.⁶ In stage 2, the agent receives the wage w and chooses an effort

⁵ We consider only these cases where subjects reciprocate based solely on effort or outcome. In reality, it is likely that the reciprocity motive is some combination of the two. The implication is that the actual outcome should be somewhere in the middle of the two proposed outcomes.

⁶ We chose the range between 0 and 14 for effort and desired effort for two reasons. First, it ensures that the maximum cost of effort is less than the maximum possible wage (cost of effort of 14 is equal to 98). Second, we wanted the efficient effort (10) to be an internal point (between 0 and 14), so that agents were not anchored to the efficient effort artificially (which could happen if the effort range was between 0 and 10).

level e that does not have to be equal to the desired effort \bar{e} specified by the contract. The cost of effort $c(e)$ is an increasing and convex function of effort, where $c(e) = e^2/2$. In stage 3, the principal first observes the outcome $y = e + \varepsilon$, which is a function of effort e and a uniformly distributed random component ε (an integer number between -2 and $+2$). As we explain below, the primary difference between treatments is what the principal can observe ($\{y, e, \varepsilon\}$ or just y). After observing y , the principal chooses an adjustment level a (an integer number between -5 and $+5$) that can be either in a form of a bonus ($a > 0$) or punishment ($a < 0$).⁷ The payoff of the principal is $\pi^P = 10y - w - |a|$ and the payoff of the agent is $\pi^A = w - c(e) + 10a$.⁸ The range of payoffs in any one period can vary substantially for both players, ranging from -105 to 160 for the principal and -148 to 150 for the agent.⁹

We employ three treatments, which we name based on what the principal observes. In the baseline Effort-Only treatment there is no random component (i.e., $\varepsilon = 0$), and the principal directly observes effort e (there is no difference between effort and outcome, since $y = e$). This treatment is similar to the bonus treatment in Fehr et al. (1997) and provides a baseline to which we compare our results. In the Effort-Shock treatment, there is a random shock component ε that the principal observes; i.e., the principal directly observes y , e , and ε . Finally, the Outcome-Only treatment is the same as the Effort-Shock treatment, but the principal only observes outcome y and does not know the composition of y .¹⁰

In all treatments, the money-maximizing subgame perfect equilibrium is for the agent to choose an effort of 0 (i.e., $e = 0$) and for the principal to make an adjustment of 0 (i.e., $a = 0$). The socially optimal

actions are for the agent to choose an effort of 10 (i.e., $e = 10$) and for the principal to provide an adjustment level of $+5$ (i.e., $a = 5$), providing a total welfare gain of 95 ($10 \times 10 - 50 + 50 - 5$).¹¹

We recruited subjects randomly from the student body of a mid-sized university in the United States. A total of 216 subjects were recruited from a standard campus-wide subject pool. Subjects interacted with each other anonymously over a local computer network. The experiment, which lasted an average of 45 minutes total, proceeded as follows. Upon arrival, subjects were randomly assigned to computer terminals and received instructions (see Online Appendix A, available as supplemental material at <http://dx.doi.org/10.1287/mnsc.2015.2177>) corresponding to one of the three treatments. The experiment was computerized using z-Tree (Fischbacher 2007). We ran 9 sessions (3 sessions per treatment) with 24 subjects in each session.

Within each session, subjects were split into 3 groups of 8 .¹² Within each group of 8 , 4 subjects were assigned to be principals and 4 were assigned to be agents. Subjects stayed in their role assignment throughout the entire experiment. In each session there were 10 periods of play. In each period subjects from opposite role assignments were randomly matched to form a principal-agent pair. After each period subjects were randomly rematched with someone of the opposite role assignment within their 8 -person group to form a new principal-agent pair. Each period proceeded in three stages. In the first stage, the principal chose a reward (an integer number between 0 and 100) and a desired effort (an integer number between 0 and 14) for the agent. After observing the reward and the desired effort, in the second stage, the agent chose an effort level (an integer number between 0 and 14). To determine the outcome, in the Outcome-Only and Effort-Shock treatments, the computer added to the effort a randomly selected number (an integer between -2 and $+2$).¹³ Then, depending on the treatment, the computer displayed

⁷ We chose the range between -5 and 5 for the adjustment because we wanted the ability to punish or reward to be large enough that most subjects would choose an internal point (to reduce censoring biases). We felt that this range accomplishes both of these goals while not being so large that contracts are completely based on bonuses or punishments.

⁸ In the treatments with shocks, the entire burden of the risk is placed on the principal, because the shock is realized after the agent chooses effort and the shock directly enters the principal's payoff function.

⁹ The principal can receive a payoff of up to 140 in the treatment without shocks. In the experiment, principals' single period payouts ranged from -95 to 150 and agents' single-period payouts ranged from -121 to 124 .

¹⁰ The two treatments that we introduce are novel, and to our knowledge have not been studied previously. However, some elements of our design are related to Xiao and Kunreuther (2015), who examine behavior in a two-person prisoner's dilemma game with stochastic versus certain outcomes, and to Cappelen et al. (2013), who study fairness views about risk taking with ex ante versus ex post stochastic outcomes.

¹¹ This does mean that nearly half of the welfare-maximizing contract comes from the bonus, although we are primarily concerned with the welfare implications of the effort chosen by the agent. Yet, since the bonus is costly and is chosen in the final stage of a one-shot game, there is no opportunity for the principal to receive any future reciprocity from the agent. Hence, we expect that most of the surplus gain will come from effort and not from the bonus. We indeed find this; Figure 2 indicates that in no treatment did more than 10% of principals give the highest possible adjustment of 5 .

¹² We divided the subjects into three groups per session in order to have three independent observations at the session level. This allows for the use of nonparametric tests, which we employ in §4. Subjects were not told that they were split into three groups of eight.

¹³ We allowed principals to receive a negative payout if effort plus the random number was negative.

to the principal either only the outcome (Outcome-Only); the outcome, effort, and the random number (Effort-Shock); or effort (Effort-Only). After observing the relevant information, in the third stage, the principal chose an adjustment level for the agent (an integer between -5 and $+5$).

At the end of each experiment, 1 of 10 periods were randomly selected for payment.¹⁴ The earnings in this period were exchanged at a rate of 10 francs = \$1. All subjects also received a participation fee of \$20 to cover potential losses. On average, subjects earned \$26 each (maximum \$42 and minimum \$7), which was paid anonymously and in cash.¹⁵

3. Predictions

Before proceeding to the results of the experiment, we provide intuition and predictions for how subjects might act under different experimental settings. To this end, we verbally discuss a model focusing on the reciprocity motive.¹⁶ After all, the experiment is centered on a double gift exchange, with agents gifting principals with effort and principals gifting agents with adjustments, and reciprocity is a key motivation of gift exchange.

It is not obvious what might motivate people to reciprocate in the context of our experiment. The literature provides some insights: numerous studies have shown that subjects often reciprocate based on both effort and outcome. For instance, Falk and Fischbacher (2006) provide a theory of reciprocity centered on the idea that people base reciprocity on both the intentions and consequences of an action. There are also recent studies on risk taking, redistribution, and charitable giving that show that some subjects condition their giving and reciprocity on both

the effort and luck of others (Charness and Levine 2007, Erkal et al. 2011, Cappelen et al. 2013, Gurdal et al. 2013, Rey-Biel et al. 2015). Hence, in this section, we present intuition for how the reciprocity motive might affect individual actions in the context of our experiment under two conditions: (1) principals reward agents based solely on effort, and (2) principals reward agents based on the outcome of the agent's actions (effort plus shock).¹⁷ In reality, it is likely that subjects will reciprocate based on some combination of effort and outcome, as the above cited papers suggest. Yet, focusing on the two extreme cases allows us to shed some light on which mechanism is more important in determining outcomes.

In the third stage, principals are able to show reciprocity for high (low) effort or outcome with a positive (negative) and costly adjustment. We assume that the principals' reciprocity motive is a relative one; i.e., it is a function of the agent's effort (or outcome) relative to the desired effort. If principals reward agents based solely on their effort and not the outcome, this means that the adjustment does not vary with the size of the shock, since the shock is outside the agent's control. On the other hand, if principals reward agents based on the outcome of their actions, then principals account for the fact that there is a component of the outcome that is beyond the agents' control, and the adjustment will vary with the shock. We also assume that agents show reciprocity as an increasing function of the wage that principals give them in the first stage. Specifically, we provide predictions for the two scenarios below.

HYPOTHESIS 1. *Principals reciprocate based on effort, meaning that principals set their adjustments based on the effort given by agents relative to desired effort, not the outcome.*

HYPOTHESIS 2. *Principals reciprocate based on outcome, meaning that principals set their adjustments based on the outcome (effort plus shock) relative to desired effort.*

In our experiment, three treatments are considered: Effort-Only (where there is no shock), Effort-Shock (where the principal sees the effort and the shock before choosing an adjustment), and Outcome-Only (where the principal sees the outcome, but not

¹⁴ When subjects are paid for multiple periods in a single experiment, the payment from one period may impact subjects' choices in another. According to Azrieli et al. (2015), paying for one randomly selected period is the only mechanism (under a wide array of assumptions) that mitigates this interperiod problem, which could otherwise cause some loss of control for the experimenter.

¹⁵ The fact that subjects receive a high participation fee of \$20 does not diminish the saliency of subject payments, because subjects may win or lose a substantial amount of money. In fact, in our experiment, some subjects made as much as \$42, while others made as little as \$7.

¹⁶ We drafted a formal model, and we found that the equilibrium outcomes are dependent on how reciprocity is modeled. The two extreme forms of reciprocity studied in this section, where reciprocity is either based solely on effort or on outcome, provide the same results in the formal model as those described in this section. Since we did not design the experiment to extract the shape of the reciprocity function (meaning that we cannot derive any meaningful testable predictions with respect to its shape) and the "extreme" results are straightforward to discuss verbally, we have only included a verbal discussion of the model to facilitate the reader's intuition for the results that we find in our analysis.

¹⁷ Focusing on the case where principals reward agents based on the fairness of the agents' actions (Rabin 1993, Fehr and Gächter 2000) gives qualitatively similar predictions. This is true, specifically, under the following fairness principle: if the agent anticipates a positive shock (with some probability), it is "fair" for the agent to split some of the surplus with the principal in the form of lower effort (as long as the outcome is also not lower), whereas if the agent anticipates a negative shock (with some probability), it is "fair" for the agent to make up for some of the lost surplus with extra effort (since the principal is the residual claimant of the lost surplus).

the effort or shock value, before choosing an adjustment). As we show below, the two hypotheses provide different testable predictions for the Effort-Only and Effort-Shock treatments. Hence, comparing these two treatments allows us to falsify at least one of the hypotheses. We therefore begin by discussing testable predictions for only the Effort-Only and Effort-Shock treatments.

Consider Hypothesis 1, where principals reciprocate based solely on effort (relative to desired effort), and agents know that principals reciprocate in this manner. In this case, the adjustment given by the principal should not vary as the shock varies. Hence, the existence of a shock should not affect the agent's effort in equilibrium (conditional on wage and desired effort) in the Effort-Shock treatment relative to the Effort-Only treatment, since the principal can perfectly observe the agent's effort and the shock. There should also be no difference in the principal's wage and desired effort in these two treatments, because the agent should be expected to react the same to these two choices in both treatments. This logic is summarized in Prediction 1.

PREDICTION 1. *If principals reciprocate based solely on effort (Hypothesis 1), then adjustments should not vary across shock levels (conditional on effort) in the Effort-Shock treatment, and there should be no difference in any of the subject's choices (wage, desired effort, effort, adjustment) in the Effort-Only and Effort-Shock treatments.*

Next, consider Hypothesis 2, where principals reciprocate based on the outcome. In this case, the principal's adjustment is increasing in the shock in the Effort-Shock treatment: as the shock increases, the outcome increases, which motivates greater reciprocation. When agents choose their effort levels, they know that there is a 0.2 probability that they will receive each of the shocks from the set $\{-2, -1, 0, 1, 2\}$. Since effort plays a smaller role in determining the adjustment in the Effort-Shock treatment than in the Effort-Only treatment (the shock plays a role in determining the adjustment in the former but not the latter), the net marginal return of higher effort is lower in the Effort-Shock treatment.¹⁸

¹⁸ To see this, assume that agents place equal weight on all five of the possible shocks and act according to the weighted sum of their different actions. One-fifth of the time there is a zero shock, and the adjustment given by principals (conditional on effort, desired effort, and wage) should be the same as in the Effort-Only treatment, since effort equals outcome. With 0.4 probability there is a positive shock. In this case, the agents want to give less effort than when there is a zero shock, but only to the extent that the outcome is not too small. Finally, with 0.4 probability there is a negative shock. Here, agents may want to give more than in the zero shock case in order to "make up" for the negative shock. However, the incentive to do so is partially mitigated by the fact that this involves a greater cost of effort, which is increasing at an increasing rate.

Since effort is costly and increasing in a convex manner, the incentive to "make up" for lost effort in the negative shock case is less than the incentive to reduce the effort created by a positive shock (due to the convex nature of the cost-of-effort curve), and the weighted effort is therefore lower in the Effort-Shock treatment than in the Effort-Only treatment. As a result, principals give a lower wage and ask for less desired effort in the Effort-Shock treatment, since their "gift" (i.e., wage), is less effective at inducing effort. This logic is summarized in Prediction 2.

PREDICTION 2. *If principals reciprocate based on the outcome (Hypothesis 2), then adjustments should be increasing in the shock in the Effort-Shock treatment, and the average effort, wage, and desired effort should be lower in the Effort-Shock treatment than in the Effort-Only treatment.*

These two predictions allow us to falsify either Hypothesis 1 or 2 (or both). Of course, we are also interested in how subjects behave in the Outcome-Only treatment; indeed, this is the treatment that is most similar to real-world principal-agent settings. In this treatment, principals do not observe effort, only the outcome. This gives agents the opportunity to "hide behind randomness" in stage 2 (Andreoni and Bernheim 2009, Aimone and Houser 2011), acting selfishly when they can ascribe their actions to chance.

We focus here on predictions for the Outcome-Only treatment under Hypothesis 2.¹⁹ In stage 3, principals can only see the outcome, not the effort or shock. Therefore, the principal should show greater reciprocity as the outcome increases (relative to desired effort) and the adjustment should be the same as in the Effort-Shock treatment. Given this logic, the agent's choice of effort in stage 2 is exactly the same in the Outcome-Only and Effort-Shock treatments. In both treatments, the agent does not know the shock value when choosing effort, and the only thing that matters to the principal when choosing the adjustment is the outcome. Hence, the decision-making calculus is the same for the agent in both treatments. From Prediction 2, this entails that the average effort, as well as the average wage and desired effort, are lower in the Outcome-Only treatment than in the Effort-Only treatment. This logic is summarized in Prediction 3.

¹⁹ We find evidence contrary to Prediction 1 in the following section, suggesting that principals do not reciprocate based solely on effort. Hence, we do not discuss predictions of the Outcome-Only treatment under Hypothesis 1. Solving for how subjects act under this hypothesis is not trivial, and the direction of the comparative statics (vis-à-vis other treatments) depends on the level of the choices in the other treatments.

Table 1 Summary Statistics

Treatment	Wage	Desired effort	Effort	Outcome	Adjustment	Principal's payoff	Agent's payoff	Total welfare
Effort-Only	41.14 (3.22)	8.95 (0.31)	6.40 (0.43)	6.40 (0.43)	0.14 (0.34)	20.91 (3.11)	15.71 (1.88)	36.62 (1.91)
Effort-Shock	33.45 (2.98)	7.63 (0.34)	4.69 (0.34)	4.62 (0.32)	−0.18 (0.22)	11.23 (3.46)	16.41 (2.49)	27.64 (1.90)
Outcome-Only	33.85 (2.28)	7.63 (0.25)	4.69 (0.41)	4.75 (0.38)	−0.50 (0.13)	12.04 (2.85)	17.63 (1.89)	29.67 (2.11)

Note. Standard errors in parentheses are based on nine independent observations.

PREDICTION 3. *If principals reciprocate based on the outcome (Hypothesis 2), then adjustments, average effort, wage, and desired effort should be the same in the Outcome-Only and Effort-Shock treatments, whereas the average effort, wage, and desired effort should be lower in the Outcome-Only treatment than in the Effort-Only treatment.*

4. Results

We observed 2,160 contracts in our experiment.²⁰ Table 1 provides the summary statistics across all three treatments. When performing statistical tests, we mainly use nonparametric tests to examine treatment effects. Each treatment has a total of nine independent observations (72 subjects per treatment, split into nine separate groups of 8 subjects each). When appropriate, we also estimate panel models with individual subjects representing random effects (to control for individual effects), standard errors clustered at the single rematching group level of 8 subjects (to control for possible correlation within a matching group), and an inverse period trend (to control for learning and experience). We consider the results starting with stage 3 first and work our way backward to stage 1.

4.1. Adjustment

In stage 3, principals choose an adjustment after seeing either the effort of the agent (in Effort-Only and Effort-Shock) or the outcome (in Outcome-Only and Effort-Shock) in stage 2. Figure 1 displays the average adjustment over periods by treatment, and Figure 2 displays the distribution of adjustment by treatment. Both the distribution and the average adjustment levels are very similar in all three treatments. Based on

the Wilcoxon rank-sum test there is no significant difference in the adjustment level between treatments: Effort-Only versus Outcome-Only (0.14 versus −0.50; p -value = 0.33, $n_1 = 9$, $n_2 = 9$), Effort-Only versus Effort-Shock (0.14 versus −0.18; p -value = 0.82, $n_1 = 9$, $n_2 = 9$), and Outcome-Only versus Effort-Shock (−0.50 versus −0.18; p -value = 0.60, $n_1 = 9$, $n_2 = 9$).²¹

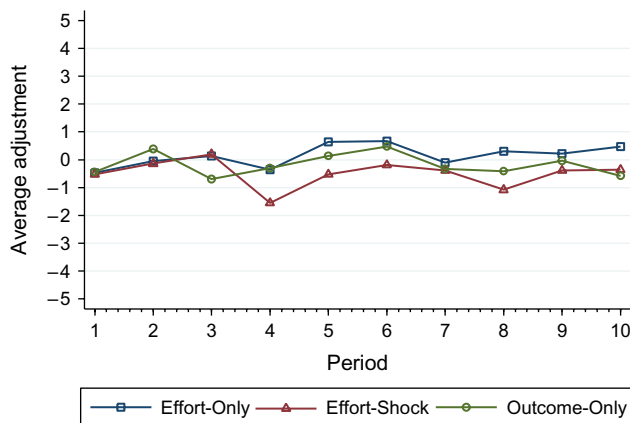
This suggests the possibility that the adjustment mechanism works relatively similarly in all three treatments. However, these results may arise from the fact that we consider the unconditional adjustment in stage 3. The model suggests that the conditional adjustment (i.e., conditional on wage, effort, desired effort, and possibly the shock) may differ across treatments, and it can also give some insight into the motivations of the subjects. Specifically, if principals reciprocate based on effort, the adjustment should not vary with the shock in the Effort-Shock treatment (Prediction 1), whereas the adjustment should vary with the shock if principals reciprocate based on outcome (Prediction 2). More generally, if the reciprocity motive is present in the principal's decision, we expect the adjustment to be a function of how "kindly" the principal was treated by the agent in stage 2. In other words, we expect the principal's adjustment to be a function of the difference between the observed effort (or outcome) in stage 2 minus the desired effort proposed in stage 1. It is also possible that the principal expects the agent to show reciprocity in stage 2 if the principal gives a large wage in stage 1, so the adjustment may also be conditional on wage.

We first test whether the effort (minus desired effort) varies across shock level in the Effort-Shock

²⁰ Of 2,160 contracts, 1,338 (62%) can be classified as individually rational and incentive compatible (IR/IC). These are contracts in which both the principal's and the agent's payoffs are nonnegative, conditional on the contract being fulfilled. Specifically, IR/IC contracts (w, e) satisfy the following two conditions: $10e - w \geq 0$ and $w - c(e) \geq 0$. We chose not to put any restrictions on the principal's decisions, because some ex ante "non-IR/IC" contracts (w, e) may be IR/IC ex post, given a certain level of adjustment a . For an experiment where the principal can only offer contracts that satisfy the IR/IC criteria, see Bartling et al. (2012). All major results hold when we focus only on the IR/IC contracts. We analyze in detail the IR/IC contracts in Online Appendix B and the non-IR/IC contracts in Online Appendix C.

²¹ We have also checked the robustness of these results using panel regression analysis. Specifically, we have estimated different panel models where individual subjects represent the random effects, and the standard errors are clustered at the single rematching group level. The dependent variable in all specifications is the *adjustment* and the independent variables are an inverse of *period*, *wage*, *effort* − *desired effort* (in Effort-Only and Effort-Shock) and *outcome* − *desired effort* (in Outcome-Only), as well as treatment dummies. All regressions indicate no significant difference in adjustment level between the three treatments. The estimation results are available in Online Appendix D.

Figure 1 (Color online) Average Adjustment by Period



treatment. Table 2 reports the average adjustment in the Effort-Only and Effort-Shock treatments as a function of whether *effort – desired effort* is negative, zero, or positive, and in the Effort-Shock treatment as a function of whether the *shock* is negative, zero, or positive. The results provide a preliminary basis for rejecting the effort-based reciprocity hypothesis, because adjustments appear to vary as the shock varies. Focusing on the case where *effort – desired effort* is negative (since *N* is high enough in this case to support statistics), we find that the average adjustment made after a negative shock is -1.11 ,

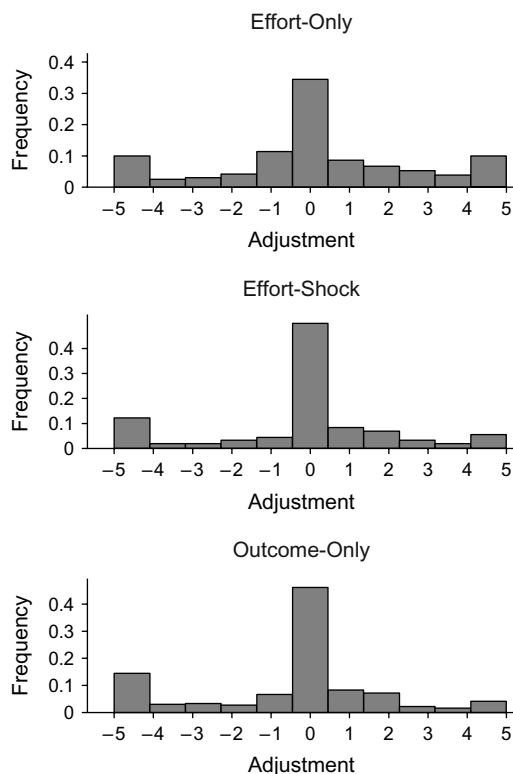
Table 2 Average Adjustment Given Effort Minus Desired Effort and Shock

	Effort – Desired effort		
	Negative	Zero	Positive
Effort-Only	-0.57	1.55	1.10
<i>N</i>	235	105	20
Effort-Shock	-0.70	0.90	1.22
<i>N</i>	252	62	46
Negative shock	-1.11	1.00	0.71
<i>N</i>	109	23	17
Zero shock	-0.55	-0.13	1.50
<i>N</i>	44	16	14
Positive shock	-0.32	1.52	1.53
<i>N</i>	99	23	15

whereas the average adjustment made after a positive shock is -0.32 . This difference is statistically significant (Mann-Whitney, p -value = 0.02). In fact, even the average adjustment made after a 0 shock, -0.55 , is marginally greater than the average adjustment made after a negative shock (p -value = 0.10).²² Since the average adjustment varies as the shock varies, we can reject the hypothesis that principals reciprocate based solely on effort.

Next, we test whether principals condition their adjustments based on previous actions. Table 3 reports the estimation results of different panel models where individual subjects represent the random effects, and standard errors are calculated using a bootstrap method.²³ The dependent variable in all specifications is the *adjustment* and the independent variables in specifications (1)–(3) are an inverse of a *period* trend, *wage*, *effort – desired effort* (in Effort-Only and Effort-Shock), and *outcome – desired effort* (in Outcome-Only).²⁴ In specifications (1) and (2), *adjustment* is positively correlated with *effort – desired effort*. In specification (3), *adjustment* is positively correlated with *outcome – desired effort*. This finding supports the idea that principals show reciprocity, since they reward higher effort (outcome) relative to desired effort.²⁵

Figure 2 Distribution of Adjustment



²² The differences in average adjustment across different shock levels are not significant when *effort – desired effort* is 0 or positive (although in two cases the p -value comes very close to significance), but this is likely due to the low number of observations in these cases.

²³ Since we have a relatively low number of clusters, we have used a bootstrap method to calculate the standard errors within each cluster (Cameron et al. 2008).

²⁴ Principals do not see effort in the Outcome-Only treatment, so we condition on outcome – desired effort.

²⁵ Table 2 suggests that principals pay a significant, and possibly discontinuous, premium for having contracts fulfilled (i.e., *effort* \geq *desired effort*). We test this by reanalyzing the results in Table 3, replacing the *effort/outcome – desired effort* variables with dummies for whether the contract was fulfilled. The coefficients on these

Table 3 Panel Models of Adjustments

Specification	(1)	(2)	(3)	(4)	(5)
Treatments	Effort-Only	Effort-Shock	Outcome-Only	Effort-Shock	Effort-Shock
<i>Wage</i>	−0.01*	−0.01**	−0.01	−0.01**	−0.01**
[<i>wage</i>]	(0.01)	(0.00)	(0.01)	(0.00)	(0.00)
<i>Effort – Desired effort</i>	0.38***	0.23***		0.23***	
[<i>effort gap</i>]	(0.08)	(0.05)		(0.05)	
<i>Outcome – Desired effort</i>			0.21***		0.23***
[<i>outcome gap</i>]			(0.06)		(0.05)
<i>Shock</i>				0.23***	
[<i>random number</i>]				(0.09)	
<i>Period</i>	−1.12***	−0.45	−0.21	−0.51	−0.51
[<i>inverse period</i>]	(0.49)	(0.37)	(0.58)	(0.39)	(0.39)
<i>Constant</i>	2.01***	0.96***	0.38	1.01***	1.00***
[<i>constant term</i>]	(0.82)	(0.33)	(0.58)	(0.33)	(0.33)
<i>N</i>	360	360	360	360	360
Clusters	9	9	9	9	9
Overall <i>R</i> -squared	0.18	0.1	0.08	0.12	0.12

Notes. The dependent variable in all specifications is the *adjustment*. Standard errors in parentheses are clustered at the group level and are calculated using a bootstrap method.

*, **, and *** represent significance at the 10%, 5%, and 1% levels, respectively.

Although these results suggest that principals reward a “kind” effort with kindness of their own, the magnitude of this reward is different across treatments. In the Effort-Only treatment, principals increase their average adjustment by 0.38 for every unit of effort given (relative to desired effort), whereas the marginal increase is only 0.23 in response to an increase in effort in the Effort-Shock treatment. In these two treatments, principals see the same information. The difference in the magnitude of these coefficients reaffirms the conjecture that principals do not reciprocate based solely on the intention (i.e., effort) of the agent. Indeed, in specification (4), we also include the *shock* as an independent variable. Consistent with our previous findings in Table 2, we find that the *adjustment* and the *shock* variables are positively correlated, suggesting that principals punish or reward agents based in part on *outcomes*. In fact, a comparison of specifications (3) and (5) suggests that principals respond similarly (0.21 versus 0.23) to an increase in *outcome* regardless of whether or not the effort is observed.²⁶

dummies are highly significant in all five specifications, and their magnitudes range from 1.61 to 2.06. The coefficient on the shock variable in column (4) remains highly significant, and none of the results reported above are qualitatively altered.

²⁶ It is possible that the relationship between “reciprocity” (adjustment) and “kindness” (effort gap or outcome gap) is not linear (Baumeister et al. 2001, Offerman 2002, Andreoni et al. 2003, Charness 2004, Bellemare and Kroger 2007). Bellemare and Kroger (2007), for example, suggest that reciprocity is a concave function of kindness (i.e., increasing in the degree of kindness increases reciprocity, but at a diminishing rate). Moreover, following the papers by Abbink et al. (2000), Fehr and Gächter (2000) and Baumeister et al. (2001), many studies have shown that “negative”

RESULT 1. There is no significant difference in the unconditional *adjustment* level between treatments. The adjustment level varies positively with effort and the shock.

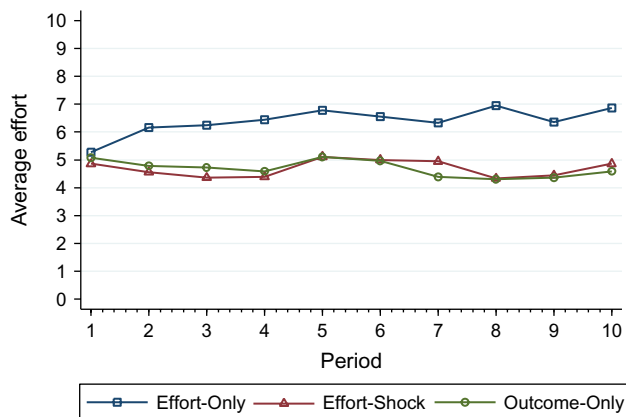
4.2. Effort

We next consider the effort that the agent chooses in stage 2. Figure 3 displays the average effort over periods by treatment, and Figure 4 displays the distribution of effort by treatment. Based on the Wilcoxon rank-sum test, we find that the average effort in the Effort-Only treatment is higher than in the Outcome-Only treatment (6.40 versus 4.69; p -value = 0.02, $n_1 = 9$, $n_2 = 9$) and the Effort-Shock treatment (6.40 versus 4.69; p -value = 0.01, $n_1 = 9$, $n_2 = 9$). On the other hand, the average effort is not different between the Outcome-Only and Effort-Shock treatments (4.69 versus 4.69; p -value = 0.82, $n_1 = 9$, $n_2 = 9$).²⁷ These results

reciprocity is stronger than “positive” reciprocity (see Charness and Kuhn 2011). In results that are available upon request, we control for both nonlinearities and distinctions between positive and negative kindness. Our results indicate that both positive and negative reciprocity increase in the degree of “kindness” (i.e., *effort – desired effort* and *outcome – desired effort* are positively correlated with *adjustment*). Moreover, positive reciprocity increases at a diminishing rate (i.e., $(\text{effort} - \text{desired effort})^2$ and $(\text{outcome} - \text{desired effort})^2$ are negative) and negative reciprocity decreases at a diminishing rate (i.e., $(\text{effort} - \text{desired effort})^2$ and $(\text{outcome} - \text{desired effort})^2$ are positive). These results are not always statistically significant, however.

²⁷ We have also checked the robustness of these results using a panel regression analysis controlling for individual subject effects, rematching groups, learning, *wage*, and *desired effort*. The regression results corroborate our main findings: effort is greater in Effort-Only than in the other two treatments, but there is no difference between Effort-Shock and Outcome-Only treatments. The estimation results are available in Online Appendix D.

Figure 3 (Color online) Average Effort by Period



are consistent with Predictions 2 and 3 of the model, which suggest that the effort given in the Effort-Only treatments is higher than in the other two treatments if principals reciprocate based on the outcome of the game play.

The intuition laid out previously indicates two reasons that agents may choose effort greater than the money-maximizing Nash prediction of zero. First, they may believe that a higher effort will lead to a greater reward (or smaller punishment) in stage 3. We showed in the previous section that such beliefs are accurate, although there are treatment differences. Second, they may exhibit positive reciprocity if the

Table 4 Panel Models of Effort

Specification	(1)	(2)	(3)
Treatments	Effort-Only	Effort-Shock	Outcome-Only
<i>Wage</i>	0.09***	0.07***	0.06***
[<i>wage</i>]	(0.01)	(0.01)	(0.01)
<i>Desired effort</i>	0.11***	0.08	0.05
[<i>desired effort</i>]	(0.05)	(0.06)	(0.05)
<i>Period</i>	−0.91**	0.39	0.09
[<i>inverse period</i>]	(0.43)	(0.64)	(0.50)
<i>Constant</i>	1.93***	1.52***	2.21***
[<i>constant term</i>]	(0.60)	(0.44)	(0.42)
<i>N</i>	360	360	360
<i>Clusters</i>	9	9	9
Overall R-squared	0.43	0.25	0.18

Notes. The dependent variable in all specifications is the subject's *effort*. Standard errors in parentheses are clustered at the group level and are calculated using a bootstrap method.

** and *** indicate statistical significance at the 5% and 1% levels, respectively.

principal gives them a high wage in the first stage; i.e., their effort is in part conditional on actions taken in stage 1. We test this possibility by conducting a panel analysis within each treatment. Table 4 reports the estimation results of different panel models, where the dependent variable in all specifications is the subject's *effort* and the independent variables are an inverse of a *period* trend, *wage*, and *desired effort*. In all specifications, there is a positive and significant relationship between *wage* and *effort*, suggesting a gift-exchange story between the principal and the agent.

RESULT 2. There is a greater *effort* in the Effort-Only treatment than in the other two treatments, whereas there is no significant difference in effort between the Effort-Shock and Outcome-Only treatments. The effort level responds positively to wage in all three treatments.

It is reasonable to suspect that the principal's willingness to reciprocate is not just a function of the absolute level of *effort* (or *outcome*), but it is also a function of the difference between *effort/outcome* and *desired effort*. Indeed, the results in Table 4 indicate that the magnitude of the effect of *desired effort* on the *effort* chosen differs between treatments. In the Effort-Only treatment, principals receive 11% of each additional unit of *effort* they desire (and this is statistically significant), whereas the magnitude is 8% in the Effort-Shock treatment and 5% in the Outcome-Only treatment (although neither are statistically significant). These results suggest that agents form reasonably correct beliefs regarding how principals will act in the adjustment period. Table 3 suggests that the adjustment response to *effort* – *desired effort* is strongest in the Effort-Only treatment, indicating that agents with correct beliefs should increase

Figure 4 Distribution of Effort

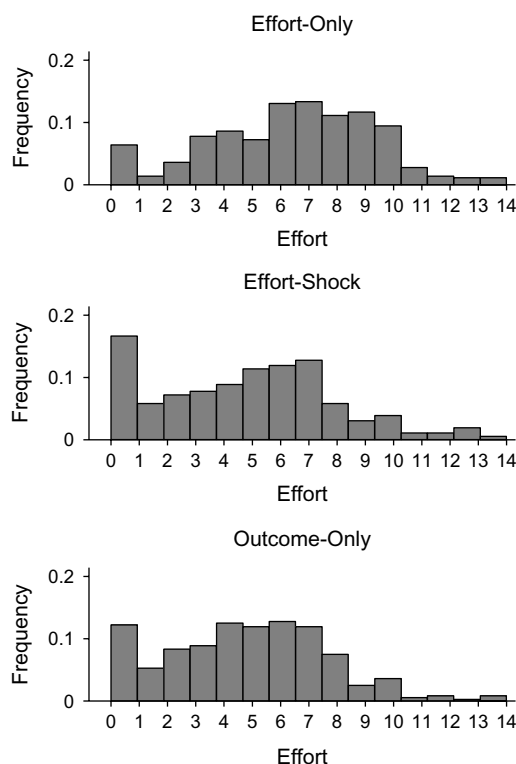
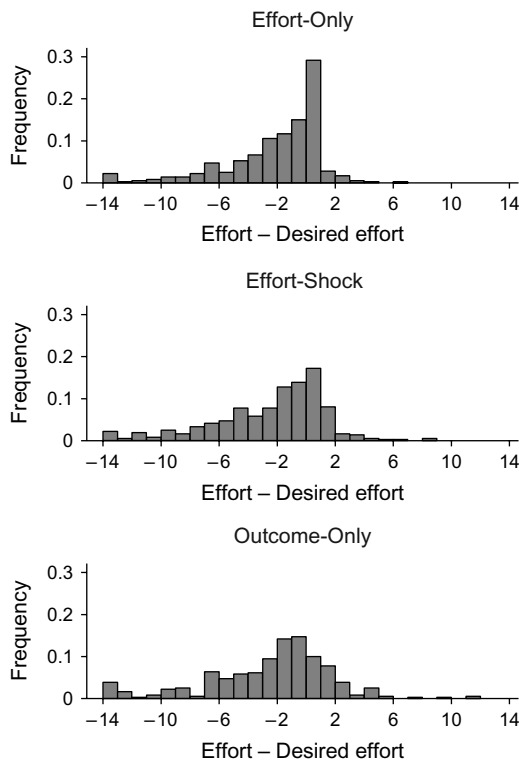


Figure 5 Distribution of Effort – Desired Effort



their effort the most in this treatment in response to an increase in desired effort.

Yet, even if agents correctly predict how principals act in stage 3, it is not clear *ex ante* how the introduction of shocks affects effort relative to desired effort. First, if agents believe that the most important thing to principals is whether the contract was fulfilled (i.e., $\text{effort}/\text{outcome} \geq \text{desired effort}$) rather than by how much it was fulfilled, we should expect to see the vast majority of effort within the interval $[-2, 2]$ of desired effort. Any effort lower than this range allows the principal to know with 100% probability that the agent did not fulfill the contract, and any effort higher than this range involves more costly effort without affecting the principal's perceived probability that the contract was fulfilled. This is precisely what we find in Figure 5, which shows the distribution of *effort – desired effort* in all three treatments. This figure indicates that the vast majority of observations in all three treatments fall within the interval $[-2, +2]$, suggesting that agents do not perceive the desired effort simply as “cheap talk” but rather as a concrete indication of the principal's expectations.

It is also quite clear from Figure 5 that the distribution of *effort – desired effort* is different in the three treatments. What can explain this? If agents are risk averse, they may choose to give more effort than desired effort in the Outcome-Only treatment (relative to the Effort-Only treatment) in order to

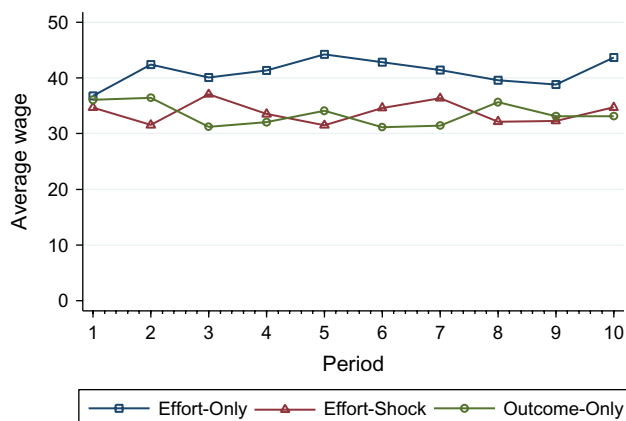
avoid any chance of being perceived as underperforming the desired effort. Whether agents choose effort greater than desired effort more frequently in the Effort-Shock treatment relative to the Effort-Only treatment depends on the degree to which agents believe that principals reward/punish based on *outcome* relative to *effort*. If agents believe that principals reciprocate based solely on effort relative to desired effort (with a discontinuity at equality), they should not give effort above the desired effort more frequently in the Effort-Shock treatment than in the Effort-Only treatment. They should give such extra effort in the Effort-Shock treatment, however, if outcomes are the primary driver of reciprocity, since negative shocks are possible. We do indeed find that the probability of effort exceeding the desired effort (based on the Wilcoxon rank-sum test) is significantly lower in the Effort-Only treatment than in the Outcome-Only treatment (0.06 versus 0.17; p -value = 0.01, $n_1 = 9$, $n_2 = 9$), is marginally lower in the Effort-Only than in the Effort-Shock treatment (0.06 versus 0.13; p -value = 0.10, $n_1 = 9$, $n_2 = 9$), and is not significantly different between the Outcome-Only and Effort-Shock treatments (0.17 versus 0.13; p -value = 0.53, $n_1 = 9$, $n_2 = 9$). This further corroborates the finding that agents expect principals to exhibit outcome-based reciprocity.

Moreover, if agents at all suspect that principals base their adjustments on outcome rather than effort, we should expect to see contracts being exactly fulfilled (i.e., effort is equal to desired effort) more often in the Effort-Only treatment than in the other two treatments. There are no shocks in this treatment, so effort is equal to outcome, whereas shocks distort the mapping from effort to outcome in the other two treatments. Our results confirm this intuition. Agents choose efforts exactly specified by the contract in the Effort-Only treatment significantly more often than in the Outcome-Only treatment (0.29 versus 0.10; p -value < 0.01, $n_1 = 9$, $n_2 = 9$) and the Effort-Shock treatment (0.29 versus 0.17; p -value = 0.03, $n_1 = 9$, $n_2 = 9$). There is no statistically significant difference between the Outcome-Only and Effort-Shock treatments (0.10 versus 0.17; p -value = 0.15, $n_1 = 9$, $n_2 = 9$).

RESULT 3. Effort levels respond positively to *desired effort* in the Effort-Only treatment. There is a greater probability that effort exceeds desired effort in the Outcome-Only and Effort-Shock treatments than in the Effort-Only treatment, whereas there is a greater probability that the contract is exactly fulfilled in the Effort-Only treatment than in the other two treatments.

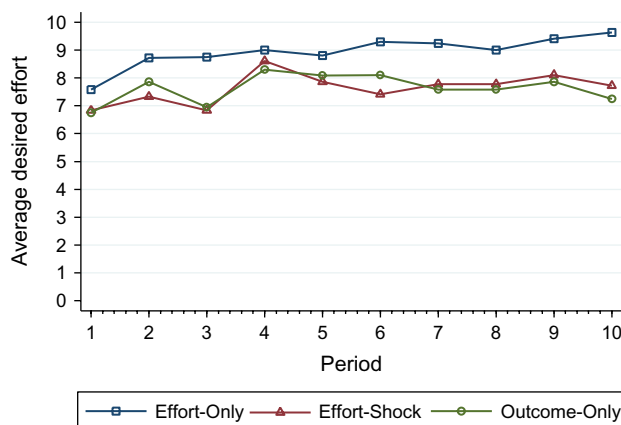
4.3. Wage and Desired Effort

In terms of welfare, the most important result presented thus far is Result 2, which indicates that

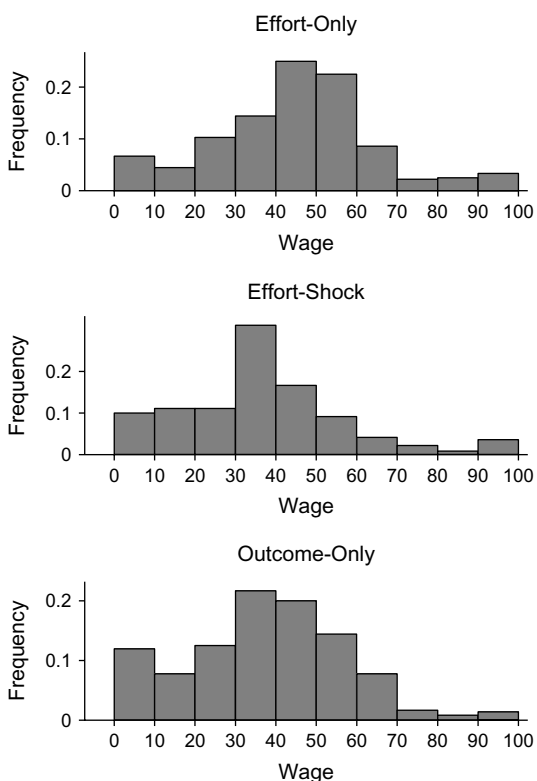
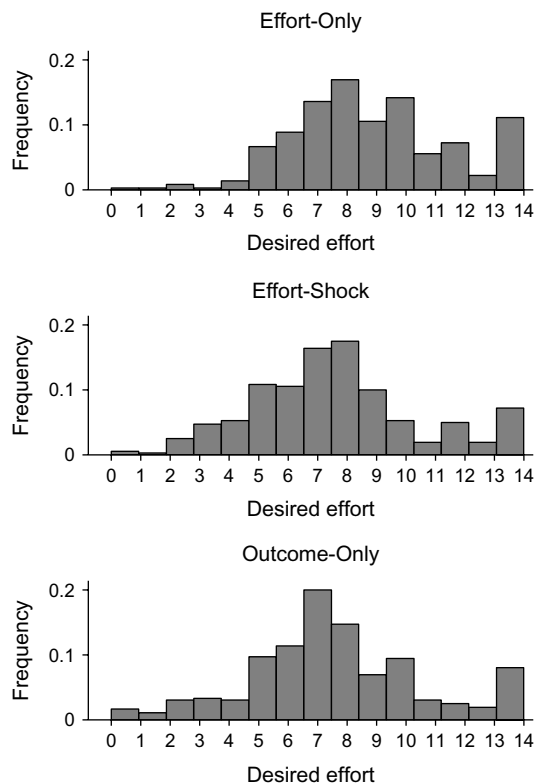
Figure 6 (Color online) Average Wage by Period

effort is greater in the Effort-Only treatment than in the other two treatments. Where does this extra effort come from? Result 3 indicates that it does not come from agents giving extra effort relative to desired effort, although it may come from agents giving less effort than desired less frequently. This leaves two non-mutually-exclusive possibilities: principals (i) give higher wages and/or (ii) ask for higher desired efforts in stage 1.

Figures 6 and 8 display the average wage and desired effort over periods and by treatment, and Figures 7 and 9 display the distribution of wage and

Figure 8 (Color online) Average Desired Effort by Period

desired effort by treatment. The average wage and desired effort are the highest in the Effort-Only treatment. Using the average within a single rematching group over all periods as one independent observation, the Wilcoxon rank-sum test shows that the average wage in the Effort-Only treatment is significantly higher than in the Effort-Shock treatment (41.14 versus 33.45; p -value = 0.05, $n_1 = 9$, $n_2 = 9$) and Outcome-Only treatment (41.14 versus 33.85; p -value = 0.08, $n_1 = 9$, $n_2 = 9$). Similarly, we find that the average desired effort is higher in the Effort-Only treatment than in the Effort-Shock treatment (8.95 ver-

Figure 7 Distribution of Wage**Figure 9** Distribution of Desired Effort

sus 7.63; p -value = 0.02, $n_1 = 9$, $n_2 = 9$) and the Outcome-Only treatment (8.95 versus 7.63; p -value = 0.01, $n_1 = 9$, $n_2 = 9$). On the other hand, wage and desired effort are not different between the Outcome-Only and Effort-Shock treatments (p -values = 0.57 and 0.89, respectively).²⁸

RESULT 4. There is a greater *wage* and *desired effort* in the Effort-Only treatment than in the other two treatments, although there is no statistically significant difference in *wage* and *desired effort* between the Effort-Shock and Outcome-Only treatments.

Result 4 indicates that the higher effort level observed in the Effort-Only treatment in Result 2 results in part from both higher wages and higher desired effort levels in the Effort-Only treatment. Why do principals offer a higher wage and ask for greater desired effort in the Effort-Only treatment? Part of the answer follows from Table 4, which indicates that effort responds significantly to desired effort in the Effort-Only treatment and not in the other two treatments. Hence, principals have more to gain from asking for higher desired effort in the Effort-Only treatment than in the other two treatments. If higher wages are necessary to induce such effort, this would also explain why *wage* is greater in the Effort-Only treatment. In fact, we find that there is a strong correlation between *wage* and *desired effort*, $\rho = 0.54$, indicating that higher wages are associated with higher desired effort.

4.4. Payoffs and Welfare

As a consequence of higher wage and higher effort, the Effort-Only treatment generates a significantly higher payoff to the principal than the Outcome-Only treatment (20.91 versus 11.23; p -value = 0.05, $n_1 = 9$, $n_2 = 9$) and the Effort-Shock treatment (20.91 versus 12.04; p -value = 0.05, $n_1 = 9$, $n_2 = 9$). Yet, the principal's payoff is not significantly different between the Outcome-Only and Effort-Shock treatments (11.23 versus 12.04; p -value = 0.89, $n_1 = 9$, $n_2 = 9$). When comparing payoffs of agents, we find no significant differences between the three treatments (all p -values > 0.48).

RESULT 5. *Principals' payoffs* in the Effort-Only treatment are higher than in the other two treatments, whereas there is no significant difference between the Effort-Shock and Outcome-Only treatments. There is no statistically significant difference in the *agents' payoffs* between any of the treatments.

²⁸ We have also checked the robustness of these results using panel regression analysis, controlling for individual subject effects, rematching groups, and learning. The regression results corroborate our main findings: wage and desired effort are greater in Effort-Only than in the other two treatments, but there is no difference between Effort-Shock and Outcome-Only treatments. The estimation results are available in Online Appendix D.

The fact that principals are better off in the Effort-Only treatment but agents are not suggests that, although principals offer higher wages in the Effort-Only treatment, this translates into higher effort levels, which leave the agents equally well off but make principals better off. The principals are made better off by enough in the Effort-Only treatment that the overall welfare (principal's payoff + agent's payoff) is greater in the Effort-Only treatment than in the other two treatments: Effort-Only versus Outcome-Only (36.62 versus 29.67; p -value = 0.05, $n_1 = 9$, $n_2 = 9$) and Effort-Only versus Effort-Shock (36.62 versus 27.64; p -value = 0.01, $n_1 = 9$, $n_2 = 9$). On the other hand, there is no significant difference in the total welfare between the Effort-Shock and Outcome-Only treatments (27.64 versus 29.67; p -value = 0.31, $n_1 = 9$, $n_2 = 9$).

RESULT 6. Total welfare is greater in the Effort-Only treatment than in the other two treatments, whereas there is no statistically significant difference between the Effort-Shock and Outcome-Only treatments.

5. Discussion and Conclusion

We conduct a gift-exchange experiment in which the agent's outcome depends on both effort and luck. Consistent with the previous literature on gift exchange (Fehr et al. 1997, 2007; Charness and Kuhn 2011), we find that bonus contracts without a shock component encourage effort and wages well above the money-maximizing Nash equilibrium prediction. We also find that a significant number of agents do not shirk and exert at least as much effort as is specified by the contract.

Two fundamental findings follow from our results. The first finding is that people reward in part on the basis of the outcome of the exchange, even if part of the outcome is determined by forces outside the control of the other party. This is not a new result. For instance, Falk and Fischbacher (2006) provide a theory of reciprocity centered on the idea that people base reciprocity on both the intentions and consequences of an action. Likewise, our result is consistent with a large literature on retrospective voting that finds voters reward/punish politicians based on outcomes over which politicians have no control (Healy et al. 2010, Gasper and Reeves 2011). It is also consistent with a large literature in psychology on outcome bias (Baron and Hershey 1988, Marshall and Mowen 1993, Mazzocco et al. 2004).

The novel and important result of our study is that the introduction of a shock in the principal-agent setting significantly reduces wages and effort, regardless of whether the shock can be observed by

the principal.²⁹ The introduction of shocks in the principal-agent setting also significantly reduces the probability of fulfilling the contract by the agent and the payoff of the principal, as well as the total welfare. The fact that shocks, even perfectly observable, have such a significant and perhaps unexpected effect in principal-agent settings has important implications for the design of optimal contracts.

What is it about the addition of shocks—observed or unobserved—that encourages principals to offer contracts with lower wages and desired effort levels? Why does the addition of shocks make agents less responsive to desired effort? Although we cannot pinpoint the exact behavioral mechanism underlying our results, we can say something about theories that are inconsistent with our results. In particular, a satisfactory theory must account for the fact that the observability of the shock does not affect effort or welfare. This suggests that our results are not being driven by agents “hiding behind randomness,” where they give less effort when they can blame a bad outcome on a negative shock (even if the shock ended up not being negative). For instance, Andreoni and Bernheim (2009) argue that people like to be perceived as fair and thus act selfishly when they can ascribe their actions to chance.³⁰ But this motivation cannot account for the multitude of differences we see between the Effort-Only and Effort-Shock treatments; since agents cannot “hide behind randomness” in the latter treatment, they should not act differently than when there is no randomness. Indeed, any explanation that cannot account for differences based on the observability of actions cannot explain our treatment differences.

What, then, can explain our results? First, our results are consistent with principals exhibiting outcome-based reciprocity. However, this simply means that we cannot reject this motivation as driving actions; it is possible that other motivations are at work as well. To this end, we believe that there are two other non-mutually-exclusive conjectures that are consistent with our results. The first conjecture has to do with the nature of gift exchange. Specifically, as wage and desired effort levels increase, the downside

risk becomes greater due to the gift-exchange nature of the game: the agent may not choose the desired effort level, and thus the higher wage is wasted. Likewise, when agents choose higher efforts, the downside risk that the principal will not reciprocate in the third stage is greater, since the effort chosen is more costly both in absolute and marginal terms. As the costs increase, players must be compensated by either higher payouts or lower uncertainty. The Effort-Only treatment offers the lowest uncertainty of the three, since agents know that principals receive an amount corresponding exactly to the amount of effort that they give. In this treatment, agents do not have to be concerned about whether the principal rewards based on intention or outcome. This, in turn, allows higher levels of effort to be sustained, because the additional risk inherent in the other two treatments makes high levels of effort too costly to be worth the risk.

Second, it may be the case that the factors affecting expected reciprocity (e.g., fairness) interact with shocks in complex ways. For example, agents may be afraid that they will be treated unfairly if they receive a bad shock in the Outcome-Only or Effort-Shock treatments. If they believe that they will be unjustly punished if they choose effort equal to the desired effort but receive a negative random number, they may instead choose effort levels lower than the desired effort, since high effort is costly. In fact, this may even be an optimal strategy in the presence of shocks. When an agent chooses effort within two levels of desired effort, the marginal gain of an additional unit of effort is only a 20% increase in the principal’s perceived probability that at least the desired effort level was given. Thus, agents have incentive at high effort levels to scale back their effort; this saves on rather large costs while minimally decreasing the probability of being perceived as choosing at least the desired effort. This effect is exacerbated if agents are averse to what they view as “unjust” punishment, since the marginal benefit to choosing at least the desired effort is lower when shocks are present.

Neither of these possibilities is mutually exclusive. In fact, they both call for further research on just how and why shocks affect contract choice. Although we know that formulating a complete, first-best contract is often not possible when shocks are present, our results suggest that the reciprocity motive does not completely mitigate this problem. Reciprocity does allow for more efficient results than standard contract theory would have us believe, but its effect is partially mitigated by the presence of shocks, whether or not the shocks are observed.

Supplemental Material

Supplemental material to this paper is available at <http://dx.doi.org/10.1287/mnsc.2015.2177>.

²⁹ Our findings contrast with the findings of Sloof and van Praag (2010), who document that subjects exert higher efforts when there is more noise in the production process. However, our results are not directly comparable since we examine behavior of subjects in a chosen-effort principal-agent setting, whereas Sloof and van Praag (2010) examine behavior in a real-effort experiment without a principal.

³⁰ Aimone and Houser (2011) also show that the “betrayal aversion” impulse is weaker when agents can hide behind randomness. In their experiment, betrayal aversion induces greater trust and hence greater efficiency. As noted above, however, this cannot explain why we do not find differences between our Effort-Only and Effort-Shock treatments.

Acknowledgments

The authors thank three anonymous referees and the editor of this journal for their valuable suggestions. They have benefitted from the helpful comments of Gary Charness, Brice Corgnet, Ron Harstad, Roberto Hernán-González, Alex Imas, Charlie Plott, David Rojo-Arjona, John List, three anonymous referees, and seminar participants at Chapman University and Case Western Reserve University, as well as participants at the North-American Economic Science Association Meeting, the 2013 International Foundation for Research in Experimental Economics Conference at Chapman University, and the 2014 Allied Social Science Associations (ASSA) Annual Meeting. The authors also thank the Economic Science Institute and Chapman University for providing facilities and Koch Foundation for financial support. Any remaining errors are the authors' own.

References

- Abbink K, Irlenbusch B, Renner E (2000) The moonlighting game: An experimental study on reciprocity and retribution. *J. Econom. Behav. Organ.* 42(2):265–277.
- Aimone JA, Houser D (2011) Beneficial betrayal aversion. *PLoS ONE* 6(3):e17725, doi:10.1371/journal.pone.0017725.
- Andreoni J, Bernheim BD (2009) Social image and the 50–50 norm: A theoretical and experimental analysis of audience effects. *Econometrica* 77(5):1607–1636.
- Andreoni J, Harbaugh W, Vesterlund L (2003) The carrot or the stick: Rewards, punishments, and cooperation. *Amer. Econom. Rev.* 93(3):893–902.
- Azrieli Y, Chambers CP, Healy PJ (2015) Incentives in Experiments: A Theoretical Analysis. Working paper, The Ohio State University, Columbus, OH.
- Baker GP (1992) Incentive contracts and performance measurement. *J. Political Econom.* 100(3):598–614.
- Baron J, Hershey JC (1988) Outcome bias in decision evaluation. *J. Personality Soc. Psych.* 54(4):569–579.
- Bartling B, Fehr E, Schmidt KM (2012) Screening, competition, and job design: Economic origins of good jobs. *Amer. Econom. Rev.* 102(2):834–864.
- Baumeister R, Bratslavsky E, Finkenauer C, Vohs K (2001) Bad is stronger than good. *Rev. General Psych.* 5(4):323–370.
- Bellemare C, Kroger S (2007) On representative social capital. *Eur. Econom. Rev.* 51(1):183–202.
- Cameron AC, Gelbach JB, Miller DL (2008) Bootstrap-based improvements for inference with clustered errors. *Rev. Econom. Statist.* 90(3):414–427.
- Cappelen AW, Konow J, Sørensen EØ, Tungodden B (2013) Just luck: An experimental study of risk-taking and fairness. *Amer. Econom. Rev.* 103(4):1398–1413.
- Charness G (2004) Attribution and reciprocity in an experimental labor market. *J. Labor Econom.* 22(3):665–688.
- Charness G, Dufwenberg M (2006) Promises and partnership. *Econometrica* 74(6):1579–1601.
- Charness G, Gneezy U (2008) What's in a name? Anonymity and social distance in dictator and ultimatum games. *J. Econom. Behav. Organ.* 68(1):29–35.
- Charness G, Haruvy E (2002) Altruism, equity, and reciprocity in a gift-exchange experiment: an encompassing approach. *Games Econom. Behav.* 40(2):203–231.
- Charness G, Kuhn P (2011) Lab labor: What can labor economists learn from the lab? Ashenfelter O, Card D, eds. *Handbook of Labor Economics* (Elsevier B.V., Amsterdam), 229–330.
- Charness G, Levine DI (2007) Intention and stochastic outcomes: An experimental study. *Econom. J.* 117(522):1051–1072.
- Ericsson KA, Charness N (1994) Expert performance: Its structure and acquisition. *Amer. Psychologist* 49(8):725–747.
- Erkal N, Gangadharan L, Nikiforakis N (2011) Relative earnings and giving in a real-effort experiment. *Amer. Econom. Rev.* 101(7):3330–3348.
- Falk A, Fischbacher U (2006) A theory of reciprocity. *Games Econom. Behav.* 54(2):293–315.
- Falk A, Kosfeld M (2006) The hidden costs of control. *Amer. Econom. Rev.* 96(5):1611–1630.
- Fehr E, Gächter S (2000) Fairness and retaliation: The economics of reciprocity. *J. Econom. Perspectives* 14(3):159–181.
- Fehr E, Schmidt KM (2007) Adding a stick to the carrot? The interaction of bonuses and fines. *Amer. Econom. Rev.* 97(2):177–181.
- Fehr E, Gächter S, Kirchsteiger G (1997) Reciprocity as a contract enforcement device: Experimental evidence. *Econometrica* 65(4):833–860.
- Fehr E, Klein A, Schmidt KM (2007) Fairness and contract design. *Econometrica* 75(1):121–154.
- Fischbacher U (2007) z-Tree: Zurich toolbox for ready-made economic experiments. *Experiment. Econom.* 10(2):171–178.
- Gasper JT, Reeves A (2011) Make it rain? Retrospection and the attentive electorate in the context of natural disasters. *Amer. J. Political Sci.* 55(2):340–355.
- Gneezy U, List JA (2006) Putting behavioral economics to work: Testing for gift exchange in labor markets using field experiments. *Econometrica* 74(5):1365–1384.
- Grossman SJ, Hart OD (1983) An analysis of the principal-agent problem. *Econometrica* 51(1):7–45.
- Gurdal M, Miller J, Rustichini A (2013) Why blame? *J. Political Econom.* 121(6):1205–1247.
- Harris M, Raviv A (1979) Optimal incentive contracts with imperfect information. *J. Econom. Theory* 20(2):231–259.
- Healy A, Malhotra N, Mo CH (2010) Irrelevant events affect voters' evaluations of government performance. *Proc. Natl. Acad. Sci. USA* 107(29):12804–12809.
- Holmström B (1979) Moral hazard and observability. *Bell J. Econom.* 10(1):74–91.
- Holmström B, Milgrom P (1991) Multitask principal-agent analyses: Incentive contracts, asset ownership, and job design. *J. Law, Econom., Organ.* 7:24–52.
- Houser D, Xiao E, McCabe K, Smith V (2008) When punishment fails: Research on sanctions, intentions, and non-cooperation. *Games Econom. Behav.* 62(2):509–532.
- Konow J (2000) Fair shares: Accountability and cognitive dissonance in allocation decisions. *Amer. Econom. Rev.* 90(4):1072–1091.
- Konow J (2003) Which is the fairest one of all? A positive analysis of justice theories. *J. Econom. Literature* 41(4):1188–1239.
- Laffont J-J, Martimort D (2002) *The Theory of Incentives: The Principal-Agent Model* (Princeton University Press, Princeton, NJ).
- List JA (2007) On the interpretation of giving in dictator games. *J. Political Econom.* 115(3):482–493.
- Marshall GW, Mowen JC (1993) An experimental investigation of the outcome bias in salesperson performance evaluations. *J. Personal Selling Sales Management* 13(3):31–47.
- Mazzocco PJ, Alické MD, Davis TL (2004) On the robustness of outcome bias: No constraint by prior culpability. *Basic Appl. Soc. Psych.* 26(2–3):131–146.
- Milgrom P, Roberts J (1992) *Economics, Organization and Management* (Prentice Hall, Englewood Cliffs, NJ).
- Offerman T (2002) Hurting hurts more than helping helps. *Eur. Econom. Rev.* 46(8):1423–1437.
- Prendergast C (1999) The provision of incentives in firms. *J. Econom. Literature* 37(1):7–63.
- Rabin M (1993) Incorporating fairness into game theory and economics. *Amer. Econom. Rev.* 83(5):1281–1302.
- Rey-Biel P, Sheremeta R, Uler N (2015) When income depends on performance and luck: The effects of culture and information on giving. ESI Working Paper, Chapman University, Orange, CA.
- Shavell S (1979) Risk sharing and incentives in the principal and agent relationship. *Bell J. Econom.* 10(1):55–73.
- Sloof R, van Praag CM (2010) The effect of noise in a performance measure on work motivation: A real effort laboratory experiment. *Labour Econom.* 17(5):751–765.
- Xiao E, Kunreuther H (2015) Punishment and cooperation in stochastic social dilemmas. *J. Conflict Resolution*, ePub ahead of print January 5, doi: 10.1177/0022002714564426.