



Dishonesty invites costly third-party punishment. Evolution and Human Behavior

Ohtsubo, Yohsuke

Masuda, Fumiko

Watanabe, Esuka

Masuchi, Ayumi

(Citation)

Evolution and Human Behavior, 31(4):259-264

(Issue Date)

2010-07

(Resource Type)

journal article

(Version)

Accepted Manuscript

(URL)

<https://hdl.handle.net/20.500.14094/90001771>



Running Head: DISHONESTY AND PUNISHMENT

Dishonesty Invites Costly Third-Party Punishment

Yohsuke Ohtsubo Fumiko Masuda Esuka Watanabe

(Kobe University)

Ayumi Masuchi

(Hokkai-Gakuen University)

This article has been accepted for publication in *Evolution and Human Behavior*.

January, 2010

Acknowledgement

This research was supported by the Japan Society for the Promotion of Science (No. 18730398). We are grateful to Tetsu Ichizawa, Keiko Koyama and Eriko Matsumoto for their help in recruiting participants, Asami Matsumura for her help in conducting the experiment, Kazuya Ishibashi for his help in preparing the figures, and Pat Barclay, Larry Fiddick, Steven Gaulin, and two anonymous reviewers for their valuable comments on earlier drafts.

Abstract

Third-party punishment for norm violators is an evolvable enforcer of social norms. The present study, involving two experiments, examined whether violations of honesty norms would induce costly third-party punishments. In both experiments, participants in the third-party role observed a protocol of the trust game, in which the trustee solicited the trustor to transfer his/her endowment by stating that the trustee would return x units from the total resource. Dishonesty was defined such that the trustee in fact returned fewer than x units. Participants were asked about their willingness to incur some cost to reduce the trustee's payoff. In Experiment 1, x was exactly half of the total resource. Participants were willing to incur more cost to punish the dishonest trustee than the trustee who allocated the resource unequally but had not sent the dishonest message. In Experiment 2, x was more than half of the total resource, and the dishonest trustee allocated the total resource equally. Therefore, the dishonest trustee was not unfair in Experiment 2. Approximately half of the participants (16 of 30) punished the dishonest but fair trustee, while few participants (1 of 30) punished the fair trustee who had not sent the dishonest message. These experiments together demonstrated that participants were willing to incur some cost to punish honesty-norm violators, even when the participants themselves were not harmed by the norm violation.

Keywords: Honesty; Social Norms; Third-Party Punishment; Trust Game

Dishonesty Invites Costly Third-Party Punishment

1. Introduction

Social norms are important foundations of human societies. Without sanctions against norm violators, however, opportunists might refrain from complying with social norms because following social norms is often personally costly (Sethi & Somanathan, 2005). Accordingly, the notion of altruistic punishments is often invoked to explain norm enforcement (Fehr & Fischbacher, 2004). The presence of sanctions may be also important for defining “social norms.” Bendor and Swistak (2001), for example, define social norms as behavioral rules involving sanctions implemented by *third parties*, arguing that the involvement of third parties distinguishes social norms from behavioral rules maintained merely by interested parties. Sanctions against a transgressor by a harmed party could be caused by non-normative motives, such as a retaliatory impulse. Therefore, the presence of third-party punishment is a yardstick whereby we can judge whether a certain behavioral rule qualifies as a social norm.

Fehr and Fischbacher (2004) experimentally demonstrated that violations of egalitarian distribution norms and cooperation norms would reliably induce third-party punishments. In the third-party punishment games, neutral third-party participants observed other players behaving in an unfair (or uncooperative) manner, and then decided to incur some cost to punish the unfair (or uncooperative) player. In Fehr and Fischbacher’s experiments, approximately 60% of the participants punished unfair (or uncooperative) players. With the same third-party punishment game, it has been shown that violations of egalitarian distribution norms induce third-party punishments across a wide range of cultures (Henrich et al., 2006). The present study applies the third-party punishment game to another, seemingly universal behavioral rule. As a traditional maxim says, “Honesty is the best policy.” Honesty norms seem to qualify as social norms.

Although none of the previous studies directly demonstrated the presence of the

third-party punishment for dishonesty, there is some suggestive evidence. In Brandts and Charness's (2003) study, for example, dyads of participants played a 2×2 game. One player (sender) was allowed to send a message indicating an intended play. Then, both players simultaneously chose their move. A sender was considered dishonest if the sender indicated that he/she would play favorably for the other player (receiver), but in fact played unfavorably for the receiver. Brandts and Charness revealed that the receivers punished the dishonest senders more severely than those who had behaved unfavorably without the dishonest message (see also Bolese, Croson, & Murnighan, 2000; Croson, Boles, & Murnighan, 2003; Pillutla & Murnighan, 1996; Schweitzer, Hershey, & Bradlow, 2006; Wang, Galinsky, & Murnighan, 2009). As we have noted, it is possible that the receivers (i.e., the "second-parties") punished the dishonest senders due to retaliatory motives. A somewhat different study conducted by Tyler, Feldman, and Reichert (2006) revealed that people disliked those who behaved dishonestly to someone else. Nonetheless, Tyler et al. did not examine whether third-party participants would engage in some costly form of punishment.

The present study aimed at testing whether dishonesty would induce costly third-party punishments. We conducted two experiments, both of which included Fehr and Fischbacher's (2004) third-party punishment game modified for the present purpose. In both experiments, focal participants observed two other participants playing the trust game, in which the first player (henceforth referred to as the "trustor") decided whether to transfer his/her endowment, 500 Japanese yen (JPY), to the second player, the trustee. If the trustor decided not to transfer the 500 JPY, both the trustor and the trustee ended the game with their initial endowment of 500 JPY. If the trustor decided to transfer the 500 JPY, the transferred endowment was tripled by the experimenter. The trustee then had 2000 JPY (i.e., his/her initial endowment of 500 JPY plus the transferred 1500 JPY) and divided it between the trustor and himself/herself. In addition to these

standard procedures, the trustee was allowed to send a pre-play message to the trustor about how the trustee would divide the money, 2000 JPY. This message, however, did not bind the trustee's later allocation behavior. The trustee is considered to have behaved dishonestly if he/she stated that he/she would return x JPY in the message but returned less than x JPY. The two experiments tested whether participants in the third-party role would be willing to incur some personal cost to punish the dishonest trustee.

2. Experiment 1

2.1. Method

2.1.1. Participants. Participants were 81 undergraduates at a middle-sized university in the *Hokkaido* area of Japan. Unbeknownst to the participants, procedurally similar experiments that included the trust game but not third-party punishment were conducted in parallel with the present study. We shall not report those studies here because their results are irrelevant to the present purpose. The 81 participants involved in the present study were randomly assigned to one of three roles (i.e., trustor, trustee, and third party). In presenting the results, we shall focus on the third-party participants. Therefore, only 27 participants (18 males and 9 females; mean age = 21.33 years, range = 20–23 with the exception of one 38-year-old participant) were included in the subsequent analyses. There were two sessions for this experiment, and each session included 11 and 16 triads, respectively.

2.1.2. Procedure. After participants arrived in a large classroom, the experimenter randomly distributed cards colored red, green, or yellow, each corresponding to one of the three roles. Participants were asked to sit according to their cards' color so that participants in the different roles sat apart. Participants were randomly assigned an ID number, whereby participants in the three roles were anonymously matched.

The experimenter handed each participant a booklet. There were three versions of the booklet, each of which described the rules of the third-party punishment game from the perspective of the participant's assigned role. The booklet explicitly explained that the trustee's message would not bind the trustee's behavior. The reward for participants was also explained as follows: After the experiment, the experimenter would randomly choose some ID numbers. Participants with those ID numbers would be paid what they earned in the game. Other participants would not be paid. The likelihood of earning a real reward was kept ambiguous for the participants.

Experiment 1 employed the strategy method, and the game was played in the following manner. First, the participants playing the trustee role were asked to indicate whether they would send a pre-determined message to their partner (i.e., the trustor). The message read: "I will give 1000 JPY back to you if you transfer your endowment." Second, the participants playing the trustor role were asked to indicate whether they would transfer their endowment in two possible cases: (i) when they received the message and (ii) when they did not. Third, the participants playing the trustee role were asked to choose one allocation scheme from the four possible schemes: (1000-1000), (1300-700), (1500-500), and (1700-300); here the left entry represents the trustee's own share, and the right entry represents the trustor's share. It was explained that this decision would be effective only when the trustor had transferred his/her endowment. Finally, the participants playing the third-party role (i.e., the focal participants) were asked whether they would punish the trustee in each of the eight possible cases: 2 (message: sent vs. not sent) \times 4 (allocation scheme). Therefore, the present experiment employed a 2×4 within-participant factorial design. When asking about their willingness to punish the trustee, we used neutral wording: e.g., "How much are you willing to pay to *reduce* the *allocator's* reward?" (italics added). It was explained to participants that they would receive 1000 JPY if their ID number was

chosen after the experiment. If they decided to use some money, c , to punish the trustee, c JPY would be subtracted from their own 1000 JPY, and $3c$ JPY from the trustee's reward. Because the third-party participants decided whether to punish the trustee for the eight cases, the third-party participants' monetary rewards would be determined depending on how the trustee behaved in the game.

After the third-party participants completed their decision task, the experimenter chose one ID number in each session and paid the monetary rewards to the chosen participants. One of the authors (AM) gave feedback about the study to the participants later in the class.

2.2. Results and Discussion

2.2.1. Results. The average amount of money that participants were willing to pay to punish the trustee is summarized in Fig. 1. This variable was submitted to a 2 (message) \times 4 (allocation scheme) repeated measures ANOVA. The main effects of message and allocation scheme were significant, $F_{1,26} = 11.33$, $P = .002$, $\eta_p^2 = .30$ for the message; $F_{3,78} = 24.22$, $P < .001$, $\eta_p^2 = .48$ for the allocation scheme. The interaction effect was also significant, $F_{3,78} = 6.46$, $P = .001$, $\eta_p^2 = .20$. The main effect of allocation scheme reflects a tendency that punishments became more severe as the chosen allocation became less fair (i.e., departed more from the equal allocation). This interpretation is confirmed by the significant linear trend associated with allocation scheme, $F_{1,26} = 3.19$, $P < .001$, (cf. quadratic and cubic trends were not significant). The main effect of message reflects the tendency that punishments were more severe when the message had been sent than when it had not. Simple effect tests indicated that there was an exception to this tendency: the effect of message was not significant in the (1000-1000) condition, $F < 1$ (see the leftmost bars in Fig. 1). Most participants did not punish the trustees who had made an equal allocation regardless of whether they had sent the message or not. The significant interaction effect was due to this exception. In sum, Experiment 1 showed that third-party punishment was

inflicted when the allocation was unfair, and that the punishment became more severe when the trustees behaved dishonestly (designated by “D” in Fig. 1) than when they did not. We also conducted a comparable ANOVA that included participants’ sex as a between-participants factor; neither main nor two-way interaction effects involving sex reached the significance level.

2.2.2. Limitations of Experiment 1. There were several limitations in Experiment 1. First, participants might not have been serious because they knew that the reward would not be paid to all participants. Second, since participants were asked to indicate their willingness to punish the trustee in all of the eight cases, participants might have attempted somehow to differentiate the two fair cases from the six unfair cases, and the five dishonesty-absent cases from the three dishonesty-present cases. In other words, the observed results might simply reflect demand characteristics due to the within-participant design (i.e., strategy method). In addition, the strategy method might diminish the role of emotions in economic games (Casari & Cason, 2009), although punitive sentiment is considered as a crucial trigger for third-party punishment (e.g., Price, Cosmides, & Tooby, 2002). Third, since all the dishonesty-present cases were associated with unfair allocations, Experiment 1 does not allow us to conclude that the participants punished dishonesty *per se*. An alternative interpretation is that the presence of dishonesty fueled the punishment for unfair behaviors, but dishonesty *per se* was not the subject of punishment. Experiment 2 was conducted to overcome these shortcomings.

3. Experiment 2

3.1. Method

3.1.1. Participants and Design. Participants were 65 undergraduates at a large university in the *Kansai* area of Japan. They volunteered to take part in this experiment in exchange for an unspecified monetary reward. Data from 60 participants (40 females and 20 males) were retained for the reported analyses after five participants were omitted, who either (i) spontaneously

revealed their suspicion about the absence of the other players during the debriefing session or (ii) erroneously believed that they could take some money from the trustee and give it to the trustor.

Experiment 2 in fact consisted of two separate experiments, which were conducted during different semesters. Although the experiments shared identical procedures and included 30 participants recruited from comparable participant pools, the experiments differed in the efficacy of punishment: the participant's payment of c JPY resulted in $3c$ JPY reduction from the trustee's reward in one experiment (henceforth Experiment 2a), and $2c$ JPY reduction in the other (Experiment 2b). Experiment 2b was conducted to test the generalizeability of the results in Experiment 2a. We shall use "Experiment 2" to refer to both Experiments 2a and 2b simultaneously.

3.1.2. Procedure. In Experiment 2, every participant was tested individually. Upon arrival, a participant was ushered to a separate room, told the nature of the experiment, and asked to sign the informed consent form. To avoid the problem of demand characteristics (i.e., not to give the participant an impression that his/her task would be to spend some of his/her money to punish another player), the experimenter told the participant that his/her reward, 1000 JPY, would be paid for his/her participation in another experiment that would require him/her to complete a lengthy questionnaire. The experimenter showed a thick questionnaire to the participant and said that he/she would be asked to complete this questionnaire after the first task. The questionnaire was, in fact, from an unrelated study, and the participant completed the questionnaire after the present study.

The explanation of the trust game was presented on a computer screen, and the participant was allowed to read through the explanation at his/her own pace. In Experiment 2, unlike in Experiment 1, the trustor and the trustee were not present. Participants were instead led

to believe that the other players were playing the trust game in other rooms. In Experiment 2, it was explained to participants that the trustee was allowed to write whatever message he/she would like.

After understanding the above instructions, the participant was presented with the trustee's message and the summary results of the trust game. The experimenter first showed the participant the message purportedly written by the trustee. In the dishonest trustee condition, the message read: "If you transfer your endowment, I will take 700 JPY and give you 1300 JPY." In the honest trustee condition, "I will allocate 1000 JPY to each" replaced the apodosis. In both conditions, the message was followed by a solicitation: "So, please transfer it to me." After the participant read the message, the experimenter showed the summary sheet of the trust game. The summary sheet indicated that the trustor transferred his/her endowment, and that the trustee allocated the endowment of 2000 JPY evenly between the two players. Notice that in both conditions, unfair allocation was not included. Therefore, punishment, if observed in Experiment 2, is not attributable to unfair behavior. The participant was then asked to indicate how much he/she was willing to pay to punish the trustee. As in Experiment 1, we used neutral wording in the instructions and questionnaires as much as possible. After the participant completed this task and completed filling out the other lengthy questionnaire, he/she was debriefed, paid 1000 JPY regardless of his/her punishment decision, and dismissed.

3.1.3. Need for the Experimental Deceptions. Experimental deceptions (i.e., leading participants to believe that there were a trustee and a trustor, who were in fact not present) were involved in Experiment 2. Although we are cognizant of some drawbacks associated with experimental deceptions (Ortmann & Hertwig, 2002), we judged them necessary in the present study because we expected that only a small portion of participants would send a dishonest message. For example, six of 27 participants (22%) did so in Experiment 1. This number could have shrunk to

nearly zero, given that Experiment 2 tested the effect of a particular combination of dishonesty and allocation behavior (i.e., a message indicating an overgenerous allocation, followed by the equal allocation). Accordingly, we decided to employ the experimental deceptions in Experiment 2.

3.2. Results and Discussion

The distribution of cost paid for punishment was highly skewed (Figs. 2a and 2b). Hence, we analyzed the punishment data as a dichotomous variable indicating whether each participant used any money for punishment. In the dishonest trustee condition, 8 of 15 participants punished the dishonest trustee in Experiment 2a, and 8 of 15 did so in Experiment 2b. In the honest trustee condition, only one participant inflicted punishment; the participant spontaneously left a comment revealing his mischievous motive (i.e., “I just wanted to see how the trustee would react, though I knew I can’t”): 1 out of 15 punished the honest trustee in Experiment 2a and 0 out of 15 in Experiment 2b. The effect of dishonesty on punishment was significant in both experiments by Fisher’s exact test, $P_s = .007$ and $.001$ (one-tailed) for Experiments 2a and 2b, respectively. As in Experiment 1, sex did not have a significant effect on punishment: 5 of 10 males and 11 of 20 females in the dishonest trustee condition punished the trustee (Experiments 2a and 2b were pooled for small n size).

The distributions of the cost paid by the punishers were also suggestive of an underlying motivation. As shown in Figs. 2a and 2b, the modal cost, excluding 0 JPY, was 100 JPY and 150 JPY in Experiments 2a and 2b, respectively. In both experiments, four of the eight punishers chose those modal costs. The modal cost in each experiment resulted in the subtraction of 300 JPY from the trustee’s reward, and thus lowered the trustee’s net reward to 700 JPY, which matched with the trustee’s dishonest statement (i.e., “I will take 700 JPY”). Hence, the results suggest the participants’ motive for eliminating the dishonesty even imperfectly (the trustor

ended the game with 1000 JPY instead of 1300 JPY).

4. General Discussion

The two experiments provided evidence that dishonesty would invite third-party punishment. Experiment 2 is considered a clearer demonstration. Unlike Experiment 1, Experiment 2 provided monetary rewards to all participants. Employing the between-participants design, and introducing an irrelevant task for which their monetary rewards were paid, we attempted to minimize the problem of demand characteristics. Experiment 2 also eliminated the confounding variable. In most of the previous studies, including Experiment 1, the manipulation of dishonesty was confounded with unfair behavior. Therefore, it was not clear whether dishonesty itself triggered the observed punishment. In Experiment 2, the dishonest trustee divided 2000 JPY equally and thus was a fair person in terms of resource allocation. The control condition (i.e., the honest trustee condition) confirmed that almost no participants would punish such a fair trustee. Nonetheless, when the equal allocation was preceded by the dishonest message, approximately half of the participants punished the dishonest but fair trustee.

The present research has an implication for the recent controversy over preferences for fairness. Falk, Fehr and Fischbacher (2008) recently tested two models of fairness preferences. One is a consequentialist model that assumes that people endorse the outcome equality. The other model assumes that the intentions that gave rise to a particular consequence matter. Falk et al. found evidence supporting the second model: when unfair allocations were made outside the allocator's control (i.e., in the absence of intention), the allocators were not punished as severely as when the allocations were made by the allocators themselves (i.e., in the presence of intention). Falk et al.'s finding is complemented by the present study's finding—fair allocations may be punished if they were mediated by a deceitful intention.

The present result is also of relevance to the reliability of human language. Based on

evolutionary game analyses, Lachmann, Számadó and Bergstrom (2001) maintain that honesty of language may be sustained by socially imposed punishment for dishonesty. The present result provides some empirical support for Lachmann et al.'s argument. It is noteworthy that Lachmann et al. argue that punishments must be restricted to self-serving dishonesty. Although we did not test this prediction, it seems reasonable that people would not punish those who stated that they would give x but in fact gave more than x . In future studies, it would be worthwhile to explore what kinds of dishonesty would or would not be punished.

In sum, the present study provides support for the notion that a behavioral rule prescribing honesty in fact qualifies as a social norm. The result implies that strong reciprocity (Fehr, Fischbacher, & Gächter, 2002; Gintis, 2000) is responsible for the evolution of honest communication among humans. The proponents of strong reciprocity generally assume that people are predisposed not only to cooperate with others but also to punish non-cooperators. The parallel assumption seems to hold in the domain of communication: people are predisposed not only to avoid telling lies (Gneezy, 2005), but also to punish dishonest others (present study). The flipside of this assumption is that dishonesty becomes more frequent when greater incentives for lying exist (otherwise no punishments are required). Recent studies, in fact, have revealed that people assess the veracity of messages more cautiously when the message senders have some incentives to behave dishonestly than when they do not (e.g., Hess & Hagen, 2006; Nakanishi & Ohtsubo, 2009). It is also expected that dishonesty is more frequent among cultures whose socio-ecological environments are associated with greater incentives for lying. To curtail such incentives, the severity and/or scope of punishment for dishonesty may vary across cultures. Future research should examine the role of third-party punishment in keeping our communication reliable in light of different levels of incentives for dishonesty.

References

- Bendor, J., & Swistak, P. (2001). The evolution of norms. *American Journal of Sociology*, 106, 1493-1545.
- Boles, T. L., Croson, R. T. A., & Murnighan, J. K. (2000). Deception and retribution in repeated ultimatum bargaining. *Organizational Behavior and Human Decision Processes*, 83, 235-259, doi:10.1006/obhd.2000.2908.
- Brandts, J., & Charness, G. (2003). Truth or consequences: An experiment. *Management Science*, 49, 116-130, doi:10.1287/mnsc.49.1.116.12755.
- Casari, M., & Cason, T. N. (2009). The strategy method lowers measured trustworthy behavior. *Economics Letters*, 103, 157-159, doi:10.1016/j.econlet.2009.03.012.
- Croson, R., Boles, T., & Murnighan, J. K. (2003). Cheap talk in bargaining experiments: Lying and threats in ultimatum games. *Journal of Economic Behavior and Organization*, 51, 143-159, doi:10.1016/S0167-2681(02)00092-6.
- Falk, A., Fehr, E., & Fischbacher, U. (2008). Testing theories of fairness—intentions matter. *Games and Economic Behavior*, 62, 287-303, doi:10.1016/j.geb.2007.06.001.
- Fehr, E., & Fischbacher, U. (2004). Third-party punishment and social norms. *Evolution and Human Behavior*, 25, 63-87, doi:10.1016/S1090-5138(04)00005-4.
- Fehr, E., Fischbacher, U., Gächter, S. (2002). Strong reciprocity, human cooperation, and the enforcement of social norms. *Human Nature*, 13, 1-25.
- Gintis, H. (2000). Strong reciprocity and human Sociality. *Journal of Theoretical Biology*, 206, 169-179, doi:10.1006/jtbi.2000.2111.
- Henrich, J., McElreath, R., Barr, A., Ensminger, J., Barrett, C., Bolyanatz, A., Cardenas, J. C., Gurven, M., Gwako, E., Henrich, N., Lesorogol, C., Marlowe, F., Tracer, D., & Ziker, J. (2006). Costly punishment across human societies. *Science*, 312, 1767-1770,

doi:10.1126/science.1127333.

Hess, N. H., & Hagen, E. H. (2006). Psychological adaptations for assessing gossip veracity.

Human Nature, 17, 337-354.

Lachmann, M., Számadó, S., & Bergstrom, C. T. (2001). Cost and conflict in animal signals and

human language. *Proceedings of the National Academy of Sciences, U.S.A.*, 98, 13189-13194,

doi:10.1073/pnas.231216498.

Nakanishi, D., & Ohtsubo, Y. (2009). Believability of secondhand social versus ecological

information in the presence of contradictory firsthand experience. *Journal of Evolutionary*

Psychology, 7, 157-166, doi:10.1556/JEP.7.2009.2.4.

Ortmann, A., & Hertwig, R. (2002). The costs of deception: Evidence from psychology.

Experimental Economics, 5, 111-131, doi:10.1023/A:1020365204768.

Pillutla, M., & Murnighan, J. K., (1996). Unfairness, anger, and spite: emotional rejections of

ultimatum offers. *Organizational Behavior and Human Decision Processes*, 68, 208-224,

doi:10.1006/obhd.1996.0100.

Price, M. E., Cosmides, L., & Tooby, J. (2002). Punitive sentiment as an anti-free rider

psychological device. *Evolution and Human Behavior*, 23, 203-231,

doi:10.1016/S1090-5138(01)00093-9.

Schweitzer, M. E., Hershey, J. C., & Bradlow, E. T. (2006). Promises and lies: Restoring violated

trust. *Organizational Behavior and Human Decision Processes*, 101, 1-19,

doi:10.1016/j.obhdp.2006.05.005.

Sethi, R., & Somanathan, E. (2005). Norm compliance and strong reciprocity. In H. Gintis, S.

Bowles, R. Boyd, & E. Fehr (Eds.), *Moral sentiments and material interests: The*

foundations of cooperation in economic life (pp. 229-250). Cambridge, MA: MIT Press.

Tyler, J. M., Feldman, R. S., & Reichert, A. (2006). The price of deceptive behavior: Disliking

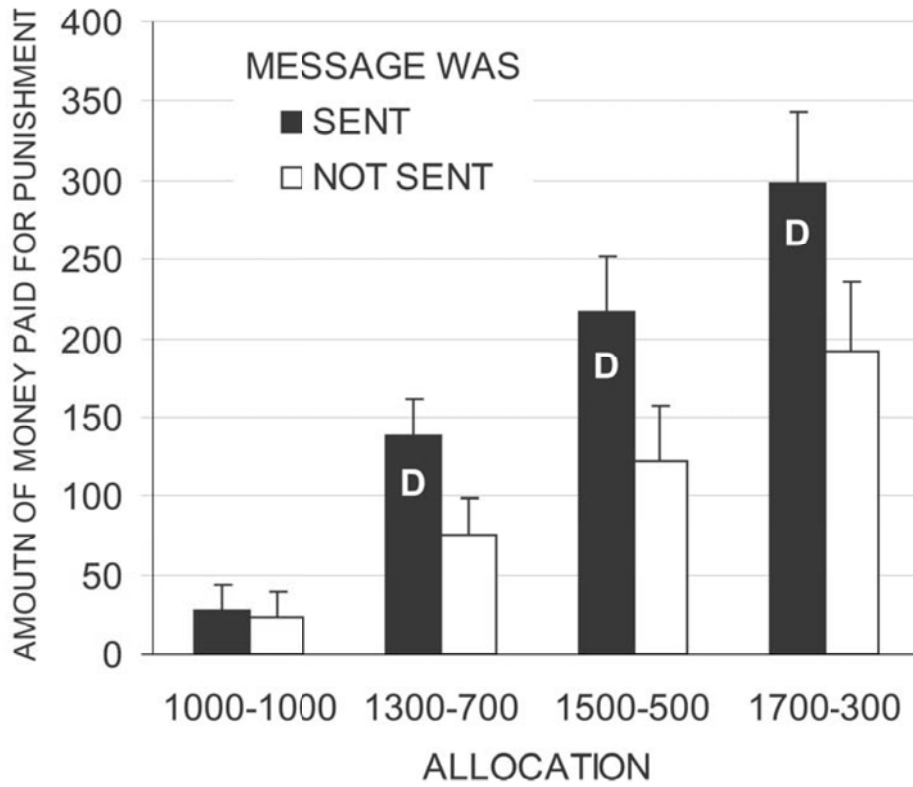
and lying to people who lie to us. *Journal of Experimental Social Psychology*, 42, 69-77, doi:10.1016/j.jesp.2005.02.003.

Wang, C. S., Galinsky, A. D., & Murnighan, J. K. (2009). Bad drives psychological reactions, but good propels behavior: Responses to honesty and deception. *Psychological Science*, 20, 634-644, doi:10.1111/j.1467-9280.2009.02344.x.

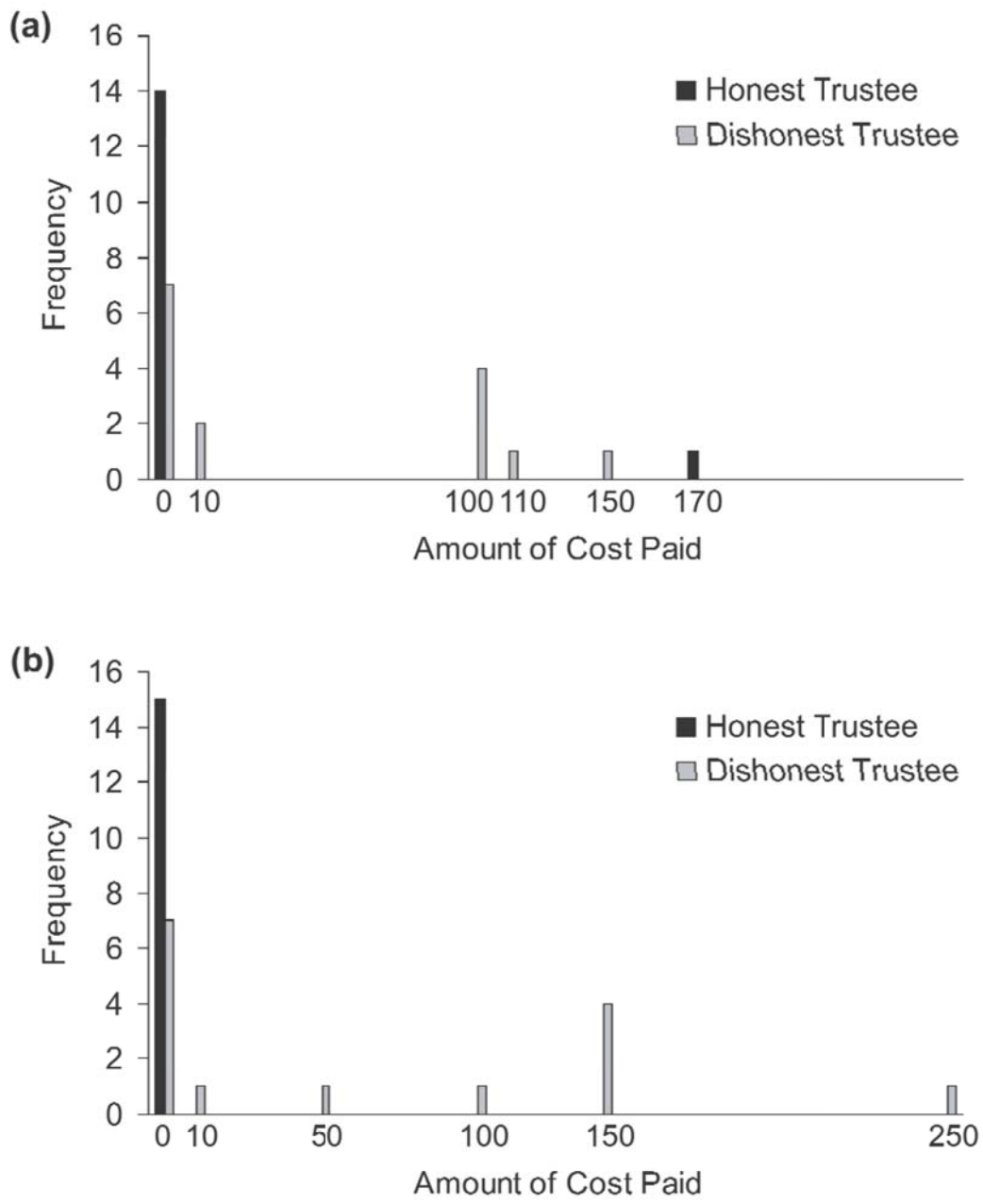
Figure Captions

Fig. 1. Mean amounts of cost paid for punishment as a function of allocation scheme and message use. “D” in bars indicates the presence of the dishonest message in the designated condition.

Fig. 2. Distributions of cost paid by participants as a function of honest/dishonest trustee condition in (a) Experiment 2a and (b) Experiment 2b.



<Figure 1>



<Figure 2>