

NBER WORKING PAPER SERIES

PUNISHMENT AND COOPERATION IN STOCHASTIC SOCIAL DILEMMAS

Erte Xiao
Howard Kunreuther

Working Paper 18458
<http://www.nber.org/papers/w18458>

NATIONAL BUREAU OF ECONOMIC RESEARCH
1050 Massachusetts Avenue
Cambridge, MA 02138
October 2012

The views expressed herein are those of the authors and do not necessarily reflect the views of the National Bureau of Economic Research.

NBER working papers are circulated for discussion and comment purposes. They have not been peer-reviewed or been subject to the review by the NBER Board of Directors that accompanies official NBER publications.

© 2012 by Erte Xiao and Howard Kunreuther. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

Punishment and Cooperation in Stochastic Social Dilemmas
Erte Xiao and Howard Kunreuther
NBER Working Paper No. 18458
October 2012
JEL No. C72,C73,C91,D02,D03

ABSTRACT

Previous findings on punishment have focused on environments in which the outcomes are known with certainty. In this paper, we conduct experiments to investigate how punishment affects cooperation in a two-person stochastic prisoner's dilemma environment where each person can decide whether or not to cooperate, and the outcomes of alternative strategies are specified probabilistically under a transparent information condition. In particular, we study two types of punishment mechanisms: 1) an unrestricted punishment mechanism: both persons can punish; and 2) a restricted punishment mechanism: only cooperators can punish non-cooperators. We show that the restricted punishment mechanism is more effective in promoting cooperative behavior than the unrestricted one in a deterministic social dilemma. More importantly, the restricted type is less effective in an environment where the outcomes are stochastic than when they are known with certainty. Our data suggest that one explanation is that non-cooperative behavior is less likely to be punished when there is outcome uncertainty. Our findings provide useful information for designing efficient incentive mechanisms to induce cooperation in a stochastic social dilemma environment.

Erte Xiao
Carnegie Mellon University
208 Porter Hall
Pittsburgh, PA, 15213
exiao@andrew.cmu.edu

Howard Kunreuther
Operations and Information Management
The Wharton School
University of Pennsylvania
3730 Walnut Street, 500 JMHH
Philadelphia, PA 19104-6366
and NBER
kunreuther@wharton.upenn.edu

I. Introduction

Punishment can be used to enforce cooperation in social dilemma situations (Yamagishi, 1986; Ostrom et al., 1992). Controlled laboratory experiments reveal that individuals are often willing to incur costs to punish defectors, even in non-repeated interactions, and that this willingness to punish can be strong enough to enforce cooperation by others (Fehr and Gächter, 2000). These studies on punishment have focused on deterministic outcome environments where agents' actions determine the outcomes with certainty. In reality, however, outcomes in social dilemmas are often determined not only by agents' actions but also by external uncertainty (Bereby-Meyer and Roth, 2006; Kunreuther et al., 2009). Such stochastic social dilemmas include a wide variety of problems involving interdependent risks in naturally-occurring environments such as airlines investing in security measures, divisions of firms undertaking risk-reducing measures to avoid a catastrophic loss that may cause the entire firm to go bankrupt, and apartment dwellers investing in fire sprinklers.

In these stochastic social dilemmas, an agent often must decide whether to incur the costs of reducing its risk of experiencing a negative outcome, knowing that even if the agent invested in a risk-reducing measure, it may still face the chance of an indirect loss if others have chosen not to follow suit. In the airline security problem, when one airline invests in baggage security, it still faces the risk of dangerous luggage transferred from other airlines unless it inspects all transferred bags (Heal and Kunreuther, 2005). An important feature of the stochastic social dilemmas is that when an agent doesn't undertake protective measures it may or may not suffer a loss due to the stochastic nature of the negative event. In fact, total payoffs will be the highest when everyone defects by choosing not to invest and if the negative outcome never occurs. We hypothesize that this feature of a stochastic social dilemma is important when considering the design of institutions to enforce cooperation, as we will elaborate below.

To understand how people punish others when there is an external outcome uncertainty, we study the impact of punishment on cooperative behavior in a two-person stochastic prisoner's dilemma game (SPD) where the outcomes of alternative strategies are specified probabilistically. In this SPD game, two agents determine whether to incur a cost to invest in protection so as to reduce the risk of losses due to the occurrence of a particular negative event. We compare behavior in this environment with actions taken in a two-person deterministic prisoner's dilemma

(DPD) game where the certain loss in the DPD game is approximately the same as the expected loss $[E(L)]$ from the negative event in the SPD game given the actions of each player.

In both games, the negative event will *not* occur if both players invest. In the DPD game, both players will suffer losses whenever one player does not invest. In the SPD game there is a well-specified probability that neither player will experience a loss even if they both decide not to invest in protection. It is important to note that the type of stochastic social dilemma situations studied in this paper differs from those situations where there is uncertainty on the information regarding what other agents have done due to imperfect monitoring (Aoyagi and Fréchet, 2009; Ambrus and Greiner, forthcoming; Grechenig et al., 2010; Patel et al., 2010; Fudenberg et al., 2012). In our experiment, there is perfect monitoring in that an agent knows what actions others have taken.¹ As we elaborate in more detail in the Literature Review (section II), the outcome uncertainty affects the punishment mechanism via different channels than behavior uncertainty.

Although numerous studies have argued that introducing peer punishment can enhance cooperation in social dilemmas (Fehr and Gächter, 2000; Ostrom et al., 1992; Yamagishi, 1988), recent research has also revealed that the unconstrained peer punishment mechanism has its limitations due to the possibility of anti-social punishment (that is, punishing cooperators) and retaliatory behavior by those who are punished.^{2,3} For example, Herrmann et al. (2008) investigate the peer punishment mechanism in public goods games using 16 different subject pools from different cities around the world such as Zurich, Seoul and Boston. Their data show that peer punishment may not be an effective means of enforcing cooperation if the cooperators are not protected from being punished. For this reason, attention has recently been focused on designing more restricted punishment mechanisms where anti-social punishment is constrained in the context of deterministic social dilemma problems (Casari and Luini, 2009; Ertan et al., 2009; Faillo et al., 2010).

¹ Some of the previous studies on imperfect monitoring also introduce noise with regard to the agents' payoff outcomes. (Aoyagi and Fréchet, 2009; Fudenberg et al., 2012)

² See Casari and Luini, 2009; Cinyabuguma et al., 2006; Denant-Boèmont et al., 2007; Dreber et al., 2008; Falk et al., 2005; Herrmann et al., 2008; Nikiforakis, 2008; Rand et al., 2010; Wu et al., 2009; Gächter et al., 2010; Gächter and Herrmann, 2011; Rand and Nowak, 2011.

³ Recent research on incentives also point out other mechanisms that may cause sanctions to backfire. For example, external incentives may crowd out intrinsic motivation (Ariely et al., 2009; Falk and Kosfeld, 2006; Fehr and Falk, 2002; Fehr and List, 2004; Fehr and Rockenbach, 2003; Frey and Oberholzer-Gee, 1997; Fuster and Meier, 2009; Gneezy and Rustichini, 2000; Herrmann et al., 2008; Houser et al., 2008; also see Gneezy et al., 2011 for a review).

In view of these previous studies on punishment, we compare the effectiveness of two punishment options in SPD and DPD supergames with a pre-specified number of periods. In Option 1, each person can incur a cost to punish her counterpart at the end of any given period after learning what strategy the counterpart pursued in that period and the resulting outcomes to both players (henceforth *BothPun*). In Option 2, only an individual investing in protection is allowed to incur a cost to punish a counterpart who has not invested in protection in that period (henceforth *InvPun*). The *InvPun* mechanism studied in this paper reflects how punishment is applied in many real world settings. For example, in formal contractual relationships when one individual reneges on his obligation, the legal system always gives the victim the right to punish the defector.

This paper contends that peer punishment mechanisms, even if they are restricted, can be less effective in promoting cooperation in the stochastic social dilemma game than in the standard (deterministic) social dilemma game because it is less clear what actions should be punished when the outcomes from one's actions involve external uncertainty. The reason is that in a stochastic environment, one's non-cooperative behavior does not necessarily lead to bad outcomes, in which case it does not impose an explicit cost on the other person. Previous studies have shown that punishment decisions are correlated with norm violations (Fehr and Gächter, 2000; Fershtman and Gneezy, 2001; Bicchieri, 2006; Xiao, 2011). In the stochastic environment, however, the outcome uncertainty may lead to normative conflict as more than one norm may exist regarding how one should behave. For example, some may think that airlines should take the risk and save the cost of investing in baggage security. Others may think that airlines should invest in baggage security to avoid the risk of a negative event. Previous research has shown that a punishment mechanism is less effective in promoting cooperation when normative conflict exists (Reuben and Riedl, 2011; Nikiforakis et al., forthcoming)

Punishment decisions are also related to the perceived negative intentions of the decision makers. In a stochastic environment, an individual may interpret non-cooperative behavior as risk-taking rather than an indication of negative intentions. The diminished perception of negative intentions can reduce punishment toward non-cooperative behavior (Blount, 1995; Nelson, 2002; Offerman, 2002; Charness and Levine, 2007; Cushman et al., 2009).

Supporting our hypothesis, we find non-investors are less likely to be punished in SPD games than DPD games. As a result, although *InvPun* mechanism is more effective than the

BothPun mechanism in promoting cooperation, it is less effective in SPD games than DPD games. Interestingly, we do not observe differences in the effectiveness of the *BothPun* mechanism in SPD games and DPD games. In particular, we find that *BothPun* does not significantly increase cooperation in either game. One explanation is that the *BothPun* mechanism leads to retaliation toward the punisher in our experiment, and less punishment means less retaliation. Therefore, less punishment does not necessarily lead to less cooperation. Our data also suggest that the occurrence of the negative event does not affect agents' punishment decisions in our setting. We discuss the policy implications of our findings at the end of the paper.

II. Literature Review

A. Stochastic social dilemma games

Uncertainty can influence the outcomes of stochastic social dilemma games in two ways: 1) when more than two players are involved (as in public goods environments), one cannot determine a specific counterpart's behavior based solely on the outcomes, or 2) agents' payoffs are affected by external risk as well as agents' actions as illustrated by the airline baggage security problem discussed above.

The effectiveness of punishment on the first type of uncertainty has recently received much attention in the economic literature with most research focusing on unrestricted peer punishment mechanisms (similar to *BothPun*). For example, Ambrus and Greiner (forthcoming) study the effect of punishment in a repeated public goods game when a subject is given information on the probability that his two counterparts took a certain action but does not know which of the two counterparts did not cooperate. They find costly peer punishment is less effective when each individual's behavior is revealed with some noise than when it is known with certainty. Grechenig et al. (2010) had a similar result. Patel et al. (2010) show that the peer punishment mechanism is less effective when there is uncertainty regarding who is a free-rider in a public goods game. Bornstein and Weisel (2010) show that punishment opportunities are not effective in promoting cooperation when information about individual endowments is incomplete in public goods games. Since the behavior of the counterpart is not known with certainty, cooperators are more likely and/or free riders are less likely to be punished than in social dilemmas without noise.

The above studies cannot inform the effectiveness of peer punishment in the second type of stochastic social dilemmas where agents' behavior is transparent but there is uncertainty related to the outcome. Furthermore, these studies did not compare the effect of uncertainty on the effectiveness of unrestricted and restricted punishment mechanisms as in our studies.

Bereby-Meyer and Roth (2006) studied games with outcome uncertainties (probabilistic PD games) and compared them with deterministic PD games. They focus on learning effects and argue that due to the noise of the payoffs, people learn to cooperate more slowly in the repeated probabilistic PD game than in the repeated deterministic PD game. Bereby-Meyer and Roth did not look at the effectiveness of punishment when there is uncertainty due to external risk. This is the focus of our paper.

B. Restricted vs. unrestricted peer punishment mechanisms and cooperation

There is a large body of research on how introducing peer punishment can enforce cooperation (Yamagishi, 1986, 1988; Ostrom et al., 1992; Fehr and Gächter, 2000; Dickinson, 2001; Andreoni et al., 2003; Fehr and Fischbacher, 2004; Fowler, 2005; Xiao and Houser, 2005, 2011; Carpenter, 2007; Sefton et al., 2007; Houser et al., 2008; also see Chaudhuri (2011) for a review).

Recently, researchers have drawn attention to the possibility of anti-social punishment and retaliation toward punishers when group members have the freedom to decide whether and who they want to punish as long as they are willing to pay the cost (Cinyabuguma et al., 2006; Dreber et al., 2008; Falk et al., 2005; Denant-Boèmont et al., 2007; Herrmann et al., 2008; Nikiforakis, 2008; Nikiforakis and Engelmann, 2011; Rand et al., 2010). Others have argued that to improve the effectiveness of punishment, it is important to restrict those who have the ability to punish (Herrmann et al., 2008). A body of research examined various forms of restricted punishment mechanisms (all in a deterministic environment) and suggests that punishment is more effective at promoting cooperation when it precludes the possibility of anti-social punishment that may emerge under the unrestricted punishment mechanism similar to the *BothPun* option. For example, Casari and Luini (2009) investigated a “consensual institution” in a public goods game whereby a request to punish a specific group member will be implemented only if at least two agents request such a punishment. Ertan et al. (2009) studied a public goods game where subjects can vote on who should be punished and found that this mechanism can

promote cooperation compared with a procedure where there were no restrictions placed on who could be punished.

Faillo et al. (2010) conducted a repeated public goods experiment where a person could punish only those who contributed less than he or she had contributed. They found that the level of cooperation doubled compared with unrestricted punishment. Casari and Plott (2003) studied a mechanism where subjects could decide whether to pay a cost to monitor others, and free-riders would be punished once discovered. They found this mechanism doubled the group's total earnings compared to when no sanction was imposed. Andreoni and Gee (forthcoming) study a "hired guy mechanism" in a public goods environment and show that people are willing to pay to introduce a delegated policing mechanism that punishes only the lowest contributor and this mechanism can avoid revenge and increase social welfare.

All these studies focus on deterministic environments and show how restricting the freedom of punishment in a certain way can improve effectiveness of the punishment mechanism. We show in this paper that when there is external uncertainty related to the outcome, a restricted punishment mechanism, *InvPun*, is less effective in promoting cooperation in an SPD setting than it is in a deterministic setting although still much more effective than the unrestricted punishment mechanism, *BothPun*.

III. Experiment

Design

As shown in Table 1, our experiment consists of six treatments determined by the baseline case and two punishment mechanisms applied to either a DPD game or an SPD game. The numbers of pairs and sessions in each treatment are given in parentheses. In each treatment, subjects are told that they will anonymously play with the same person for 10 periods, after which they will be matched with a different participant to play the same game for another 10 periods.

Baseline cases without punishment In the baseline SPD game, two subjects are paired and play together for 10 periods. At the beginning of each period, each subject is given 48 talers (2 talers = \$1). The two players must make a decision simultaneously on whether to invest in a risk-reduction measure to prevent a random negative event (with a loss of 24 talers). If both players choose to INVEST, then the negative event will not happen. The investment cost to each player

is 12 talers. If only one player invests, then there is a 40% chance that the negative event will happen. In other words, there is a 40% chance the investor will lose 36 talers and the non-investor will lose 24 talers, and there is a 60% chance that the investor will lose 12 talers and the non-investor will lose 0 talers. If both choose to NOT INVEST, then there is a 64% chance that the negative event will happen, in which case each player will lose 24 talers. Thus, there is a 36% chance that each will lose 0 talers in this situation. We chose these parameters because, based on previous results on SPD games (see Kunreuther, et. al. 2009), the cooperation level is not very high in this setting. It thus allows us to investigate how punishment mechanisms may promote cooperation.⁴

In a DPD game without punishment treatment, the known losses are the expected losses of the corresponding scenario of the SPD game.⁵ At the end of each period, each player is informed about: 1) her counterpart's decision; 2) whether the negative event occurred; and 3) total losses to each player. The payoff tables for the SPD and DPD games are shown in Tables 2(a) and 2(b), respectively.

Treatments with Punishment In the *DPD_BothPun* treatment, after each agent makes a decision on whether or not to invest, she sees the counterpart's decision and her own current payoff. She then decides whether to punish the counterpart. The *DPD_InvPun* treatment differs from the *DPD_BothPun* treatment in that only those who have invested can punish a counterpart if the counterpart has *not* invested.

In the *SPD_BothPun* treatment, agents decide whether or not to punish their counterparts after they are informed of their counterpart's decision and whether or not the negative event occurred. We reveal the outcome of the negative event to the subjects before they make the punishment decision because it allows us to examine whether the punishment decision depends on whether or not the negative event occurred. In the *SPD_InvPun* treatment, an agent who invested can, after being informed whether or not the negative event occurred, determine whether or not to punish her counterpart, but only if the counterpart has **not** invested.

⁴ For the both risk seeking and risk neutral subjects, the dominant strategy in SPD game is to not invest. For subjects who are highly risk averse, it is possible that (Invest, Invest) is a Nash Equilibrium. For example, if $U(x)=x^{1-r}$ where r is the risk aversion coefficient and x is the payoff, for (Invest, Invest) to become a Nash Equilibrium, r has to be greater than 1.1.

⁵ We use the integer numbers that are closest to the expected payoffs in SPD game to make it easier for subjects to compare the payoffs of different strategies. This also allows subjects to receive integer amount of payoffs in all the treatments.

In all the punishment treatments, subjects must determine whether to punish the other person, and if so by how much, before learning about their counterpart's punishment decision. In each period, a punisher can have up to 24 talers deducted from the counterpart's payoff. Every 3 talers deducted from the counterpart costs the punisher 1 taler. In the *BothPun* treatments, subjects might have negative earnings in one or more periods. In this case, they are told that their earnings for that period were zero. At the end of each period, each subject learns what her earnings were in that period and whether her counterpart inflicted any punishment, and if so, its severity.

Procedure

We conducted the experiment at the Behavioral Lab at the Wharton School, University of Pennsylvania by recruiting 310 subjects from the general student population at the University of Pennsylvania. Each subject participated in only one treatment. The experiment was conducted using Z-tree (Fischbacher, 2007). Subjects were told that, in addition to a fix payment of \$10, one pair would be randomly selected at the end of the experiment and would receive their actual payments from a single period that was also randomly chosen.⁶ Each subject was randomly assigned to one computer terminal. Before the experiment started, each subject completed an exercise to make sure he or she understood the payoffs under different strategies.

IV. Hypotheses

Hypothesis 1: The BothPun mechanism is less effective than the InvPun mechanism in promoting cooperation in both an DPD and SPD game.

Unlike the public goods game with punishment opportunities, where subjects are often unsure about who punished them (Fehr and Gächter, 2000; Herrmann et al., 2008), those who

⁶ We recognize that the incentive might be different when one pair of subjects is selected for a payoff as opposed to a situation where everyone is paid according to their decisions. (For more discussion see Holt, 1986 and Camerer and Hogarth, 1999). On the other hand, different numbers of subjects participated in each session and we did not observe any behavioral pattern suggesting a correlation between the decision on whether or not to cooperate and the probability of receiving payment. The fact that we do observe difference in subjects' punishment decisions that is consistent with our hypothesis also suggests that subjects do take their decisions seriously.

receive punishment in two-player PD games always know who punished them. This knowledge may lead to more retaliation than in a game where there are two or more counterparts.

Furthermore, retaliatory punishment may lead to anti-social punishment where cooperators are punished when they inflict punishment on others (Falk et al., 2005; Herrmann et al., 2008; Nikiforakis, 2008; Dreber et al., 2008; Rand et al., 2010).

In a repeated game, the punishment opportunity under the *BothPun* mechanism has a dual role: norm enforcing and retaliation. An individual is more likely to cooperate when she expects her counterpart to incur costs for punishing non-cooperative behavior (i.e., norm-enforcing impacts). However, being punished can also lead to retaliation against the punisher. This retaliatory punishment can lead to less cooperation. The *BothPun* mechanism thus leads to less cooperation than in the baseline case with no punishment if the detrimental effect of retaliatory punishment is greater than the positive norm-enforcing effect of punishment. In contrast, the *InvPun* mechanism places constraints on anti-social and retaliatory punishments. Under this mechanism, for a person to retaliate, one has to first cooperate. We thus hypothesize that *InvPun* will promote greater cooperation than the baseline treatment where no punishment is allowed, and also, greater cooperation than the *BothPun* mechanism.

Hypothesis 2: Compared to a DPD game, in a SPD game it is less clear whether a person should be punished for not cooperating (that is, not investing). Specifically, under the BothPun mechanism, the difference in the probability of being punished when one invested and when one did not invest is smaller in an SPD game than a DPD game. Under the InvPun mechanism, non-investors are less likely to be punished in the SPD game as compared to those in the DPD games.

Previous social dilemma studies with deterministic outcomes show that the non-cooperators are much more likely to be punished than cooperators (Feher and Gächter, 2002). In those studies, like the DPD game, non-cooperative decisions will reduce the payoff for the group. In an SPD game, however, non-cooperative behavior (in our experiment, to not invest) does not necessarily lead to lower earnings. In fact, both individuals earn the highest payoff if neither invests and the negative event does not occur. We hypothesize that, compared to a DPD game, this feature of an SPD game makes it less clear whether a person should be punished for not investing for the following reasons:

First, previous research has shown that normative conflict can influence the effectiveness of punishment mechanisms in enforcing cooperation (Reuben and Riedl, 2011; Nikiforakis et al., forthcoming). Normative conflict may arise in an SPD game due to outcome uncertainty. In particular, some people may approve of risk-taking behavior and view the decision to not invest as optimal whether or not other agents decided to invest in protection since agents may be spared losses even if they do not undertake these measures. Others are likely to view non-investment behavior as inappropriate since they are more likely to suffer large losses than if their counterpart had invested in protective measures. Second, the decision to punish is found to be positively correlated with perceived negative intentions of the counterpart (Blount, 1995; Nelson, 2002; Offerman, 2002; Charness and Levine, 2007). The outcome uncertainty in the SPD game clouds the intentions underlying the non-investment decisions of the counterpart in the SPD game: non-cooperative behavior might simply reflect risk preferences and agents do not know their counterpart's risk preference. In contrast, a non-cooperative decision in a DPD game clearly reveals the individual's intention.⁷

⁷ Note that while people's perception regarding appropriate behavior in SPD game may be correlated with their own risk preference (probably due to self-serving bias, see Nikiforakis et al., forthcoming), how people interpret the

More specifically, we hypothesize that under the *BothPun* mechanism in a DPD game, non-investors are more likely to be punished than investors. However, the difference in the probability of being punished when one invested and when one did not invest is smaller in the SPD game than the DPD game. Under the *InvPun* mechanism, non-investors are less likely to be punished in the SPD game as compared with those in the DPD game.

We next discuss the implication of these two hypotheses on the effectiveness of the two punishment mechanisms in SPD and DPD games.

Based on our discussion of Hypothesis 1, the difference in punishment behavior between SPD and DPD environments may have a dual impact on the effectiveness of the *BothPun* punishment mechanism in promoting cooperation. More specifically, less implementation of punishment on non-investors may diminish the norm-enforcing function of punishment but it can also lead to less retaliation, which can help maintain cooperation. Thus, the relative effectiveness of punishment in enforcing cooperation under the *BothPun* mechanism depends on the magnitude of the two effects in the SPD and DPD games.

In contrast, the *InvPun* mechanism places constraints on anti-social and retaliatory punishments. We expect that when subjects are less likely to punish their counterparts who did not invest, the norm-enforcing function of punishment will be diminished. Therefore, we predict that the *InvPun* mechanism is less effective in promoting cooperation in the SPD game than in a DPD game because a non-investor is less likely to be punished when outcomes are stochastic than when they are known with certainty.

intentions may not correlate with their own risk preference but rather their belief about the counterpart's risk preference. In this paper, we focus on the behavioral consequences of outcome uncertainty in a *repeated* interaction and provide behavioral evidence supporting our hypothesis on how outcome uncertainty affects the punishment decisions. Future studies (probably based on a *one-shot* SPD game) would be valuable to investigate the correlation between subjects' risk preferences and their perception of the norm; subjects' belief of the counterparts' risk preference and their belief regarding the intention of their counterparts' behavior.

V. Results

Each subject played the supergame with one partner for 10 periods and with a different partner for another 10 periods. Data in both 10-period sequences support our hypotheses, although we observe some differences between the two sequences. Here, we report the results from the data in the first 10 periods where subjects did not have any experience.⁸

We first report the aggregate analysis of investment rate in each treatment, then investigate how subjects make their punishment decisions and how investment decisions are affected by punishment in the previous period.

Aggregate analysis

Figure 1 plots the dynamics of investment rates over time in each treatment. It shows that in both the DPD and SPD environments, the *InvPun* mechanism achieves a significantly higher cooperation rate in almost every period than the *BothPun* mechanism. Table 3 reports the average investment rate in each treatment. Supporting Hypothesis 1, In the DPD and SPD contexts, the *BothPun* and *InvPun* mechanisms increase the average investment rate compared to the baseline treatment, but only the increase under the *InvPun* mechanism is statistically significant (see Mann-Whitney test results in Table 3).⁹

⁸ Data analysis and results of the second 10 periods are reported in an earlier version of this paper and are available on request.

⁹ Previous studies often show that unrestricted punishment can promote cooperation in public goods environments although restricted ones are more effective. The ineffectiveness of *BothPun* in promoting cooperation in our studies may be caused by the large amount of retaliatory and anti-social punishment that occurred in the experiment as we report below. Unlike a public goods game, in a two-person PD game, it is his/her counterpart who imposed the punishment. The Dreber et al. (2008) study is probably the most directly comparable to our study; they find a substantial increase in cooperation when punishment is allowed. However, their design differs from ours in that subjects can choose *either* to punish *or* defect in the next round while subjects in our study can choose both to punish and defect which may lead to a more intensive retaliatory environment.

Table 3 also shows that *DPD_InvPun* achieves a statistically significantly higher investment rate than the *SPD_InvPun* treatment (84.77% vs. 65.43%, Mann-Whitney test, $p=0.01$), although the investment rate in the corresponding baseline SPD and DPD treatments are approximately the same. We conducted a Tobit regression analysis using each pair's average investment rate over the 10 periods as the dependent variable and the six treatments' dummy variables as the independent variables. The regression result is reported in Table 4. We find the difference between the coefficients of *SPD_InvPun* and SPD significantly differs from that between the coefficients of *DPD_InvPun* and DPD (F-test, $p=0.04$). These results provide evidence that the *InvPun* mechanism is more effective in promoting cooperation in a DPD game than in an SPD game. Such differences in the cooperation rate between a DPD and a SPD game are not observed under the *BothPun* mechanism (F-test, $p=0.70$).

We also examine whether punishment mechanisms promote efficiency by comparing average earnings with the corresponding baseline treatment. The data is reported in Table 3. We find that in SPD and DPD games, the earnings are about the same when punishment is available as when it is not. In particular, the average earnings (in talers) is 33.3 in SPD and 30.9 in *SPD_BothPun*, and 33.1 in *SPD_InvPun*. None of the pairwise comparisons is significant (Mann-Whitney test, $p>0.10$). Similarly, the average earnings (in talers) in DPD is 33.5, 30.8 in *DPD_BothPun* and 33.6 in *DPD_InvPun*. None of the pairwise comparisons is significant (Mann-Whitney test, $p>0.10$). This suggests that even though the restricted punishment mechanism is more effective in promoting cooperation, the benefit from punishment implementation does not exceed its cost.

Individual analyses of punishment decisions

Descriptive analyses Table 5 reports descriptive data of punishment decisions in the SPD and DPD games. The data suggests a substantial amount of punishment behavior in all treatments. Punishment tends to be much less frequent and severe in the *InvPun* mechanism than in the *BothPun* mechanism.

Figure 2 plots the percentage of investors/non-investors who received punishment in each treatment. We separate the case when the subject punished the counterpart in the previous period from the case when she did not.¹⁰ The comparison of these two cases within each treatment indicates to what degree punishment may be triggered by retaliation. Figure 2 shows that, in the *BothPun* option, the frequency of receiving punishment is much higher when one punished the counterpart in the previous period than when one did not punish the counterpart, but this is not the case in the *InvPun* option for both the DPD and SPD games. Thus our data show that the *InvPun* mechanism significantly reduces retaliatory punishment. For example, in the SPD_BothPun treatment when subjects did not invest, they are punished 70% of the time (42 of 60 cases) if they inflicted punishment on their counterparts in the previous period. In contrast, if the non-investors did not punish their counterparts in the previous period, they are punished only 8% of the time (19 of 239 cases). This difference suggests that the implementation of punishment is dominated by retaliatory motives.

Figure 2 also suggests that, in the *BothPun* treatments, there is a substantial amount of anti-social punishment. For example, investors who punished their counterparts in period $t-1$ are punished 42% of the time (11 of 26 cases) in period t in the DPD_BothPun treatment, and 76% of the time (31 of 41 cases) in period t in the SPD_BothPun treatment.

¹⁰We also tried to analyze the data further to separate the cases when the individual invested in the previous round from when she did not. However, under such a classification, we had fewer than 10 observations in some cells. Moreover, the data does not provide much new information.

Regression analyses Statistical evidence for the presence of retaliatory punishment is provided by a random individual effect ordered probit regression analysis of punishment decisions. We find an individual is significantly more likely to punish her counterpart when she was punished in the previous period (chi-square test, $p < 0.01$, see Appendix A Table A-1 for the regression result).

One potential difference between SPD and DPD games is that non-cooperators may receive less punishment in the SPD game if individuals' decisions on whether to impose a costly punishment depend on whether they experienced a negative outcome due to *outcome bias* (see Baron and Hershey, 1988; Cushman et al., 2009). To examine this, we calculate the frequency of punishment when a loss occurred and when it did not in an SPD game. In the *SPD_BothPun* treatment, following the occurrence of a negative event, the frequency of punishment is 17% compared to 19% when the negative event did not happen. In the *SPD_InvPun* treatment, for the cases where one invested and the counterpart did not invest, about 62% investors punished the non-investor counterpart when the negative event occurred and about 53% investors did so when the negative event did not occur. We conducted a random individual effect ordered probit regression analysis of individuals' punishment decisions and found the occurrence of a loss does not have a significant effect on the punishment decisions (see Appendix A Table A-2 for the regression results).¹¹ This result suggests that an outcome bias in judgment and choices does not necessarily extend to punishment decisions in the repeated stochastic social dilemma environments.

¹¹ We also undertook non-parametric tests; the results are consistent with the regression results. For each group, we calculated the average expenses subjects paid to punish the counterpart during the 10 periods for each of the six scenarios specified in the Table A-2 regression (1) or the two scenarios (negative event occurred or not) specified in Table A-2 regression (2). We found the occurrence of the negative event does not have a significant effect on punishment decisions (Wilcoxon sign-ranked tests, $p > 0.05$).

One possible explanation is that, as we discussed Hypothesis 2, due to the uncertainty of a negative outcome in SPD game, there might be a normative conflict so that it is less clear whether an Invest or Not Invest decision should be punished in the SPD environments. That is, in the stochastic environment, people may interpret defection as risk-taking rather than a norm violation or an indication of negative intentions. Thus, even if the outcome bias leads people to think the non-cooperative decision is worse when the bad outcome occurs, it may not change their preference on punishment as it is just considered risk-taking behavior rather than being morally wrong.

To test Hypothesis 2, we compare the difference in the proportion of subjects being punished when they invested and when they did not invest in the SPD and DPD environments in the *BothPun* case. In the *InvPun* mechanism only non-investors can be punished and therefore, we compare only the fraction of non-investors who are punished in the SPD and DPD environments. About 31% of non-investors and 8% investors were punished in *DPD_BothPun*. The difference in the fraction of subjects being punished when they invested and when they did not invest is much smaller in the *SPD_BothPun* treatment (21% non-investors and 16% investors are punished). Under the *InvPun* mechanism, when the counterpart did not invest, about 86% in the DPD game and only 56% in the SPD game were punished.

To provide statistical evidence for the difference in punishment behavior between SPD and DPD games using the subject's punishment amount as the dependent variable, we conducted a random individual effect ordered probit regression analysis in *SPD_BothPun* and *DPD_BothPun*. We use three categories in order to ensure that the number of observations in each cell is sufficiently large. The dependent variable is ($CPunAmtReceived_{i,t}$) which equals: (i) "0" if subject i received no punishment; (ii) "1" if subject i received a positive punishment

amount that is no more than 12; (iii) "2" if subject i received punishment exceeding 12. We include the subject's investment decision as the independent variable ($Inv_{i,t}$ and $NotInv_{i,t}$) and allow different coefficients for each treatment ($SPD_BothPun$ or $DPD_BothPun$). Given the results of retaliatory punishment above, we also allow different coefficients for the cases when the subject punished the counterpart in the previous period ($Punisher_{i,t-1}=1$) and those where she did not ($NPunisher_{i,t-1}=1$) (see Table 6 Regression (1)).¹²

We found the difference between the coefficient of " $Inv_{i,t} * NPunisher_{i,t-1} * SPD_BothPun$ " (-1.70) and " $NotInv_{i,t} * NPunisher_{i,t-1} * SPD_BothPun$ " (-1.09) is significantly less than the difference between the coefficient of " $Inv_{i,t} * NPunisher_{i,t-1} * DPD_BothPun$ " (-2.02) and of " $NotInv_{i,t} * NPunisher_{i,t-1} * DPD_BothPun$ " (-0.42) (chi-square test, $p=0.02$).¹³ This result suggests that when subjects did not punish the counterpart (excluding the possibility of retaliatory punishment by the counterpart, $NPunisher_{i,t-1}=1$), the difference in the punishment amount received by investors ($Inv_{i,t}=1$) compared with those received by non-investors ($NotInv_{i,t}=1$) is significantly smaller in the SPD environment than in the DPD environment. This result from *BothPun* conditions provides first statistical evidence supporting Hypothesis 2.

¹² We also tried regressions that control for group effect. For example, we conducted random individual effect OLS regression analysis of punishment amount received using standard error clustering on pairs. (We ran the regressions using STATA,; this software does not allow clustering in a random individual effect ordered probit regression). We still find non-investors are less likely to be punished in SPD environments and DPD environments (chi-square tests, one tail, $p<0.01$ for the comparison between $SPD_BothPun$ and $DPD_BothPun$; $p=0.03$ for the comparison between SPD_InvPun and DPD_InvPun). We also added the counterpart's punishment and investment decisions in the previous period in the random individual effect ordered probit regression and the difference is still significant (chi-square tests, one tail, $p<0.01$).

¹³ " $Inv_{i,t} * NPunisher_{i,t-1} * SPD_BothPun$ " = 1 if subject i in $SPD_BothPun$ treatment *invested* in period t and did not punish the counterpart in the previous period; =0, o.w.

" $NotInv_{i,t} * NPunisher_{i,t-1} * SPD_BothPun$ " = 1 if subject i in $SPD_BothPun$ treatment *did not invest* in period t and did not punish the counterpart in the previous period; =0, o.w.

" $Inv_{i,t} * NPunisher_{i,t-1} * DPD_BothPun$ " = 1 if subject i in $DPD_BothPun$ treatment *invested* in period t and did not punish the counterpart in the previous period; =0, o.w.

" $NotInv_{i,t} * NPunisher_{i,t-1} * DPD_BothPun$ " = 1 if subject i in $DPD_BothPun$ treatment *did not invest* in period t and did not punish the counterpart in the previous period; =0, o.w.

We ran a similar regression analysis to compare the punishment amount received by non-investors between the *SPD_InvPun* and *DPD_InvPun* treatments (see Table 6 Regression (2)). Since in these two treatments, only non-investors can be punished and only when their counterpart invested, we include only those observations when subject i did not invest and her counterpart invested in period t . We find again that when subjects did not punish the counterpart in the previous period, non-investors are less likely to be punished in the *SPD_InvPun* than the *DPD_InvPun* treatment. More specifically, the coefficient of “ $NPunisher_{i,t-1} * SPD_InvPun$ ” is significantly smaller than that of “ $NPunisher_{i,t-1} * DPD_InvPun$ ” (0.15 vs. 1.51, chi-square test, $p < 0.01$). Thus, the results from *InvPun* treatments also support Hypothesis 2.

In summary, consistent with previous studies (Nikiforakis, 2008; Dreber et al., 2008) we find evidence of retaliating punishment in *DPD* games. In addition, we find retaliating punishment also occurs in *SPD* games. Interestingly, the occurrence of the negative event does not affect punishment decisions. More importantly, supporting our hypothesis, the punishment behavior suggests that it is less clear what behavior should be punished in the *SPD* environment than in the *DPD* environment.

Effect of punishment on investment

We ran a random individual effect probit regression analysis of each individual’s investment decision in period t for each punishment treatment. The independent variables include whether the individual i invested in period $t-1$ ($Inv_{i,t-1}$), whether the counterpart invested in period $t-1$ ($CtpInv_{i,t-1}$), and the amount of punishment that individual i received in the period $t-1$ ($PunAmtReceived_{i,t-1}$). As we pointed out earlier, punishment may have a detrimental effect when it is anti-social (i.e., punishment is imposed on the investors); therefore, we allow the

investor and the non-investor to have different coefficients for this variable ($NotInv_{i,t-1} * PunAmtReceived_{i,t-1}$ and $Inv_{i,t-1} * PunAmtReceived_{i,t-1}$) in the *BothPun* treatment for the SPD and DPD games. We first ran the regression including all the above variables, a period variable and the dummy of the end period. Since these last two variables are not significant and they do not change the conclusions detailed below, we report only the regression result of the model that does not include the period variables.¹⁴

The regression results reported in Table 7 suggests that the coefficient of “ $Inv_{i,t-1} * PunAmtReceived_{i,t-1}$ ” (i.e., the amount of punishment subject i received in the previous period when she invested) is significantly negative in both the *SPD_BothPun* and *DPD_BothPun* treatments ($p < 0.01$). This implies that under the *BothPun* mechanism, the larger the punishment incurred in period $t-1$ when a person invested, the less likely the person will invest in period t .

On the other hand, as shown in Table 7, if one did not invest in period $t-1$, the punishment amount does not have any significant positive effect on investment decisions in period t . More specifically, the coefficient of $NotInv_{i,t-1} * PunAmtReceived_{i,t-1}$ in the *SPD_BothPun* (0.01) is not significant and, in fact, it is marginally significantly negative in the *DPD_BothPun* treatment (-0.03). In contrast, the coefficient of $PunAmtReceived_{i,t-1}$ is positive in both the *SPD_InvPun* and *DPD_InvPun* treatments although it is significant only in the *SPD_InvPun* treatment ($p = 0.04$).

VI. Conclusion

¹⁴ These regression models control for group effect in that all four models include the counterpart’s investment decisions and their punishment decisions in the previous period. We also undertook random individual effect OLS regression analysis of the investment decision using standard error clustering on pairs. We find the coefficient of $Inv_{i,t-1} * PunAmtReceived_{i,t-1}$ is still significantly negative in the *DPD_BothPun* treatment ($p < 0.01$) although it is not significant in the *SPD_BothPun* treatment. Also, the coefficient of “ $PunAmtReceived_{i,t-1}$ ” is significantly positive in the *SPD_InvPun* treatment ($p = 0.04$).

We conducted experiments to investigate the impact of peer punishment on promoting cooperation in stochastic social dilemma environments where the payoffs are decided not only by agents' behavior, but also by some exogenous uncertainty. In particular, we studied two types of punishment mechanisms: an unrestricted one where both individuals can punish (*BothPun*) and a restricted one where only investors can punish non-investors (*InvPun*) and compared behavior with a baseline case of no punishment. We find the *InvPun* treatment increases cooperation relative to the baseline case and the *BothPun* treatment in the DPD and SPD games. However, the *InvPun* mechanism is less effective when there is external uncertainty related to the outcome.

Our study contributes to the understanding of how exogenous outcome uncertainty affects people's decisions in imposing and reacting to peer punishment. We provide supporting evidence for the hypothesis that non-investors are less likely to be punished in a more realistic stochastic environment than when outcomes are known with certainty. As a result, although a restricted punishment mechanism where only investors can punish non-investors is more effective for promoting cooperation than the unrestricted punishment in the DPD game, the enhanced effectiveness of punishment mechanism is less significant in the SPD game. This finding suggests that for peer punishment mechanisms to be more effective in a stochastic social dilemma environment, it may not be enough to exclude anti-social punishment. When outcomes are uncertain, it is important to resolve normative conflicts and convey clear normative messages to the community what behavior is disapproved and should be punished. Future studies are needed to explore how different instruments that help to convey normative messages can be

applied with or without punishment mechanisms to enforce cooperation in stochastic social dilemma environments.¹⁵

¹⁵ See Kunreuther, Meyer and Michel-Kerjan (in press) for more details on why individuals do not invest in protective measures and the rationale for building codes and other regulations to reduce losses from natural disasters.

References

- Aoyagi, M. and Fréchette, G. (2009). Collusion as public monitoring becomes noisy: Experimental evidence. *Journal of Economic Theory*, 144 (2009) 1135–1165.
- Ambrus, A. and Greiner, B., forthcoming. Imperfect public monitoring with costly punishment - An experimental study. *American Economic Review*.
- Andreoni, J., Harbaugh, W. and Vesterlund, L. (2003). The Carrot or the Stick: Reward Punishment and Cooperation. *American Economic Review*, 2003, 93, pp. 893-902.
- Andreoni, J. and Gee, L. (forthcoming) Gun For Hire: Delegated Enforcement and Peer Punishment in Public Goods Provision. *Journal of Public Economics*.
- Ariely, D., Bracha, A. and Meier, S. (2009). Doing Good or Doing Well? Image Motivation and Monetary Incentives in Behaving Prosocially. *American Economic Review*, 99(1), 544 - 555.
- Baron J. and Hershey, J.C. (1988). Outcome bias in decision evaluation. *Journal of Personality and Social Psychology*, 54(4) Apr, 569-579.
- Bereby-Meyer, Y. and Roth, A. (2006). The speed of learning in noisy games: Partial reinforcement and the sustainability of cooperation. *American Economic Review*, 96, 1029-1042.
- Bicchieri, C., (2006). *The Grammar of Society: The Nature and Dynamics of Social Norms*, Cambridge University Press.
- Blount, S. (1995). When Social Outcomes Aren't Fair: The Effect of Causal Attributions on Preferences. *Organizational Behavior and Human Decision Process*, 63(2), 131-144.
- Bornstein, G. and Weisel, O. (2010). Punishment, cooperation, and cheater detection in "noisy" social exchange. *Games* 1, 18-33.
- Camerer, C. F and Hogarth, R. M, (1999). "The Effects of Financial Incentives in Experiments: A Review of Capital-Labor-Production Framework," *Journal of Risk and Uncertainty*, 19(1-3), pp. 7-42.
- Carpenter, J. (2007). Punishing Free-Riders: How Group Size Affects Mutual Monitoring and the Provision of Public Goods. *Games and Economic Behavior*, 60(1): 31-52.
- Casari, M. and Luini, L. (2009). Group cooperation under alternative punishment institutions: an experiment" *Journal of Economic Behavior and Organization*, 71(2), 273-282.
- Casari, M. and Plott, C. R. (2003). Decentralized management of common property resources: experiments with a centuries-old institution. *Journal of Economic Behavior & Organization* ,51, 217-247.
- Charness, G. and Levine, D. (2007). Intention and Stochastic Outcomes: An Experimental Study. *Economic Journal*, 117, 1051-1072.
- Chaudhuri, A. (2011). Sustaining cooperation in laboratory public goods experiments: a Selective Survey of the Literature. *Experimental Economics*, 14(1), 47-83.
- Cinyabuguma, M., Page, T. and Putterman, L, (2006). Can second-order punishment deter perverse punishment? *Experiment Economics*, 9(3), 265-279.
- Cushman, F.A., Dreber, A., Wang, Y., and Costa, J. (2009). Accidental outcomes guide punishment in a 'trembling hand' game. *PLoS ONE* 4(8): e6699. doi:10.1371/journal.pone.0006699.
- Denant-Boèmont, L., Masclet, D., and Noussair, C.N. (2007). Punishment, counterpunishment, and sanction enforcement in a social dilemma experiment. *Economic Theory*, 33, 145-167.
- Dickinson, D. (2001). The Carrot vs. the Stick in Work Team Motivation.. *Experimental Economics*, 4, 107-124.

- Dreber, A., Rand, D., Fudenberg, D. and Nowak, M. (2008). Winners don't punish. *Nature*, 452, 348-351.
- Ertan, A., Page, T. and Putterman, L. (2009). Who to punish? Individual decisions and majority rule in mitigating the free rider problem. *European Economic Review*, 53, 495-511.
- Faillo, M., Grieco, D. and Zarri, L. (2010). Legitimate Punishment, Feedback, and the Enforcement of Cooperation. Working paper No. 16, Department of Economics, University of Verona.
- Falk, A., Fehr, E. and Fischbacher, U. (2005). Driving Forces Behind Informal Sanctions. *Econometrica*, 73 (6), 2017-2030
- Falk, A. and Kosfeld, M. (2006). The Hidden Cost of Control. *American Economic Review*, 96(5), 1611-1630.
- Fehr, E. and Falk, A. (2002). Psychological Foundation of Incentives. *European Economic Review*, 46, 687-724.
- Fehr, E. and Fischbacher, U. (2004). Social norms and human cooperation. *TRENDS in Cognitive Sciences*, 2004, 8(4), 85-190.
- Fehr, E. and Gächter, S. (2000). Cooperation and Punishment in Public Goods Experiments. *American Economic Review*, 90(4), 980-994.
- Fehr, E. and List, J. (2004). The Hidden Costs and Rewards of Incentives. *Journal of the European Economic Association*, 2(5), 743-771.
- Fehr, E. and Rockenbach, B. (2003). Detrimental Effects of Sanctions on Human Altruism. *Nature*, 422, 137-140.
- Fershtman, C., and Gneezy, U. (2001) Strategic Delegation: An Experiment. *RAND Journal of Economics*, Vol. 32, No.2, 352-368.
- Fischbacher, U. (2007). z-Tree: Zurich Toolbox for Ready-made Economic Experiments. *Experimental Economics*, 10(2), 171-178.
- Fowler, J.H. (2005). Altruistic punishment and the origin of cooperation. *Proceedings of National Academy of Sciences*, 102, 7027-7049.
- Frey, B. and Oberholzer-Gee, F. (1997). The Cost of Price Incentives: An Empirical Analysis of Motivation Crowding-Out. *American Economic Review*, 87, 746-755.
- Fudenberg D., Rand, D.G. and Dreber, A. (2012) Slow to Anger and Fast to Forgive: Cooperation in an Uncertain World. *American Economic Review*, 102 720-749.
- Fuster, A. and Meier, S. (2009). Another Hidden Cost of Incentives: The Detrimental Effect on Norm Enforcement. *Management Science*, doi 10.1287/mnsc.1090.1081.
- Gächter, S., Herrmann, B., and Thoeni, C. (2010). Culture and Cooperation. *Philosophical Transactions of the Royal Society B – Biological Sciences*, 365(1553), 2651-2661
- Gächter, S. and Herrmann, B. (2011). The limits of self-governance when cooperators get punished: Experimental evidence from urban and rural Russia. *European Economic Review*, 55(2), 193-210
- Gneezy, U. and Rustichini, A. (2000). A Fine Is A Price. *Journal of Legal Studies*, 29(1), 1-17.
- Gneezy, U. S. Meier and P. Rey-Biel (2011). When and Why Incentives (Don't) Work to Modify Behavior. *Journal of Economic Perspectives*, 25, 4, 191-210.
- Grechenig, C., Nicklisch, A. and Thöni, C. (2010). Punishment despite reasonable doubt - a public goods experiment with uncertainty over contributions. *Journal of Empirical Legal Studies*, 7, 847-867.
- Heal, G.M. and Kunreuther, H. (2005). IDS Models of Airline Security. *Journal of Conflict Resolution*, 41:201-17.

- Herrmann, B., Thöni, C. and Gächter, S. (2008). Antisocial Punishment Across Societies. *Science*, Vol. 319. no. 5868, pp. 1362 – 1367.
- Holt, C. (1986). “Preference reversals and the independence axiom.” *American Economic Review*. 76: 508-515.
- Houser, D., Xiao, E., McCabe, K. and Smith, V. (2008). When punishment fails: Research on sanctions, intentions and non-cooperation. *Games and Economic Behavior*. 62(2), 509-532.
- Kunreuther, H., Meyer, R.J. and Michel-Kerjan, E. (in press). Strategies for better protection against catastrophic risks. In: E. Shafir (Ed.) *Behavioral Foundations of Policy*, Princeton, NJ: Princeton University Press.
- Kunreuther, H., Silvasi, G., Bradlow, E. and Small, D. (2009). Deterministic and Stochastic Prisoner's Dilemma Games: Experiments in Interdependent Security. *Judgment and Decision Making*, 4(5), 363–384.
- Nelson W.R. Jr. (2002). Equity and intention: it's the thought that counts. *Journal of Economic Behavior and Organizations*, 48(4), 423-430.
- Nikiforakis, N. (2008). Punishment and counter-punishment in public good games: Can we really govern ourselves? *Journal of Public Economics*, 92(1-2), 91-112.
- Nikiforakis, N., Noussair, C. and Wilkening, T. forthcoming. Normative Conflict & Feuds: The Limits of Self-Enforcement. *Journal of Public Economics* (in press).
- Nikiforakis, N. and Engelmann, D. (2011) Altruistic Punishment and the Threat of Feuds. *Journal of Economic Behavior and Organization*, 78 (3), 319-332.
- Offerman, T. (2002). Hurting Hurts More Than Helping Helps. *European Economic Review*, 46, 1423-1437.
- Ostrom, E., Walker, J. and Gardner, R. (1992). Covenants With and Without a Sword: Self-Governance is Possible. *American Political Science Review*, 86, 404 – 417
- Patel, A., Cartwright, E., and Van Vugt, M. (2010) Punishment Cannot Sustain Cooperation in a Public Good Game with Free-Rider Anonymity. Göteborg University Department of Economics Working Papers in Economics 451.
- Rand, D.G., Armao, J.J., Nakamaru, M. and Ohtsuki, H. (2010). Anti-social punishment can prevent the co-evolution of punishment and cooperation. *Journal of Theoretical Biology*, 265, 624-632
- Rand, D.G. and Nowak, M.A. (2011). The evolution of antisocial punishment in optional public goods games. *Nature Communications*, 2 434.
- Reuben, E. and Riedl, A. (2011). Enforcement of Contribution Norms in Public Good Games with Heterogeneous Populations. Working paper. IZA.
- Sefton, M., Shupp, R. and Walker, J. (2007). The Effect of Rewards and Sanctions in Provision of Public Goods. *Economic Inquiry*, 45(4) 671 – 690.
- Wu, J., Zhang, B., Zhou, Z., He, Q., Cressman, R., Tao, Y. (2009). Costly punishment does not always increase cooperation. *Proceedings of the National Academy of Sciences*, vol. 106 no. 41 17448-17451
- Xiao, E. (2011). Profit seeking punishment corrupts norm obedience. Working paper, Carnegie Mellon University.
- Xiao, E. and Houser, D. (2005). Emotion Expression in Human Punishment Behavior. *Proceedings of the National Academy of Sciences*, 102(20), 7398-7401.
- Xiao, E. and Houser, D. (2011) Punish in Public. *Journal of Public Economics*, 95, 1006–1017.

- Yamagishi, T. (1986). The provision of a sanctioning system as a public good. *Journal of Personality and Social Psychology*, 51, 110-116.
- Yamagishi, T. (1988), Seriousness of Social Dilemmas and the Provision of a Sanctioning System. *Social Psychology Quarterly*, 51(1) 32-42.

Table 1. Conditions and treatments and number of pairs and sessions in each condition/treatment

Treatments	Conditions	SPD	DPD
No punishment (Baseline)		SPD (32) (6 sessions)	DPD (24) (4 sessions)
Subjects can punish each other		SPD_BothPun (31) (5 sessions)	DPD_BothPun (23) (4 sessions)
Only investors can punish non-investor counterparts		SPD_InvPun (23) (5 sessions)	DPD_InvPun (22) (4 sessions)

Note: The number in the first parenthesis is the number of pairs in each treatment.

Table 2.a. Possible outcomes in the SPD game

		Agent j	
		Invest	Not Invest
Agent i	Invest	(-12, -12)	40%: (-36, -24) 60%: (-12, 0)
	Not Invest	40%: (-24, -36) 60%: (0, -12)	64%: (-24, -24) 36%: (0, 0)

Table 2.b. Possible outcomes in the DPD game

		Agent j	
		Invest	Not Invest
Agent i	Invest	(-12, -12)	(-22, -10)
	Not Invest	(-10, -22)	(-15, -15)

Note: The left number in each cell is agent i 's loss and the right number is agent j 's loss. The loss outcomes in the DPD game in each scenario equal the expected value of the losses in the SPD game of the corresponding scenario (we use the integers that are closest to the expected payoffs in SPD game).

Table 3. Investment percentage and average earnings by treatment

Treatment	Obs.	Investment percentage (%)		Average earnings	
		Mean (s.e.)	p-value *	Mean (s.e.)	p-value *
SPD	32	36.88 (6.15)		33.30 (1.01)	
SPD_BothPun	31	48.23 (6.84)	0.21	30.89 (1.21)	0.17
SPD_InvPun	23	65.43 (6.43)	0.01	33.11 (0.60)	0.45
DPD	24	33.96 (6.82)		33.47 (0.22)	
DPD_BothPun	23	51.52 (8.44)	0.10	30.83 (1.26)	0.51
DPD_InvPun	22	84.77 (5.24)	< 0.01	33.61 (0.70)	0.10

Note: standard error is calculated using each pair as one observation.

* This column reports the p-value of Mann-Whitney test of the investment percentage (or average earnings) between the punishment treatment and the baseline without punishment treatment under each condition. We calculate the average investment rate (or average earnings) of 10 periods for each pair. We count each pair as one observation.

Table 4. Tobit regression analysis of average investment rate of each group

	Dependent variable: average investment rate over 10 periods of each group	
	Coef.	s.e.
SPD	0.36***	0.08
SPD_BothPun	0.51***	0.08
SPD_InvPun	0.69***	0.09
DPD	0.30***	0.09
DPD_BothPun	0.52***	0.09
DPD_InvPun	1.00***	0.10

*** significant at 1% level; **significant at 5% level; *significant at 10% level

Table 5. Punishment decisions by treatment

Treatment	Punishment frequency (%)	Amount of talers paid to punish			
		Mean	Median	95th percentile	Max
SPD_BothPun	18.23	0.90	0	8.00	8
SPD_InvPun	8.70	0.32	0	2.50	8
DPD_BothPun	19.13	0.84	0	7.00	8
DPD_InvPun	7.27	0.44	0	5.50	8

Table 6. Random individual effect ordered probit regression analysis of punishment amount received

	Dependent variable: <i>CPunAmtReceived</i> $_{i,t}$ (=0 if subject i' received no punishment in period t ; =1 if $0 <$ subject i's received punishment amount ≤ 12 ; =2 if subject i's received punishment amount > 12)			
	Regression (1)		Regression (2)*	
	Coef.	s.e.	Coef.	s.e.
$Inv_{i,t} * Punisher_{i,t-1} * SPD_BothPun$	-0.03	0.31		
$Inv_{i,t} * NPunisher_{i,t-1} * SPD_BothPun$	-1.70***	0.32		
$NotInv_{i,t} * NPunisher_{i,t-1} * SPD_BothPun$	-1.09***	0.29		
$Inv_{i,t} * Punisher_{i,t-1} * DPD_BothPun$	0.01	0.39		
$NotInv_{i,t} * Punisher_{i,t-1} * DPD_BothPun$	-0.02	0.37		
$Inv_{i,t} * NPunisher_{i,t-1} * DPD_BothPun$	-2.02***	0.40		
$NotInv_{i,t} * NPunisher_{i,t-1} * DPD_BothPun$	-0.42	0.33		
$NPunisher_{i,t-1} * SPD_InvPun$			0.15	0.49
$Punisher_{i,t-1} * DPD_InvPun$			0.89	0.90
$NPunisher_{i,t-1} * DPD_InvPun$			1.51**	0.67
cut1_cons	0.64	0.31	-0.19	0.47
cut2_cons	1.53	0.32	0.96	0.49

Note: *** significant at 1% level; **significant at 5% level; *significant at 10% level

$Inv_{i,t}=1$ if i invested in period t ; =0, o.w. $NotInv_{i,t}=1$ if i did not invest in period t ; =0, o.w.

$NPunisher_{i,t-1}=1$ if i did not punish the counterpart in period t ; =0, o.w. $Punisher_{i,t-1}=1$ if i punished the counterpart in period t ; =0, o.w.

The baseline for each regression is the case when subject i in SPD treatments did not invested in period i and punished the counterpart in period t-1.

*includes only the cases where the individual did not invest but her counterpart invested in the current period.

Table 7. Random individual effect probit regression analysis of investment decisions

Dependent variable: Invest _{i,t} =1 if i invest in period t; =0, o.w.				
	SPD_BothPun	SPD_InvPun	DPD_BothPun	DPD_InvPun
	Coef. (s.e.)	Coef. (s.e.)	Coef. (s.e.)	Coef. (s.e.)
Inv _{i,t-1}	1.36*** (0.24)	0.37 (0.23)	1.86*** (0.21)	1.35*** (0.42)
CtpInv _{i,t-1}	1.28*** (0.16)	0.54*** (0.20)	1.24*** (0.20)	0.85*** (0.34)
NotInv _{i,t-1} *	0.01 (0.01)		-0.03* (0.02)	
PunAmtReceived _{i,t-1}				
Inv _{i,t-1} *	-0.05*** (0.02)		-0.10*** (0.03)	
PunAmtReceived _{i,t-1}		0.05** (0.02)		0.03 (0.02)
_cons	-1.39*** (0.18)	-0.11 (0.29)	-1.46*** (0.15)	-0.13 (0.65)

Note:

Inv_{i,t-1}=1 if i invested in period t; =0, o.w.

CtpInv_{i,t-1}=1 if i's counterpart invested in period t; =0, o.w.

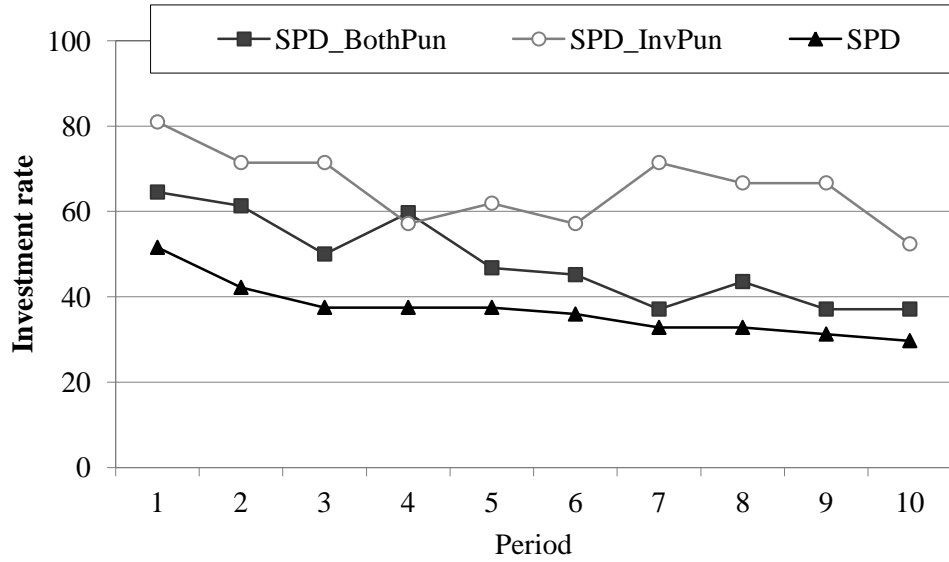
NotInv_{i,t-1}=1 if i did not invest in period t; =0, o.w.

PunAmtReceived_{i,t-1}: the punishment amount imposed on subject i by her counterpart in period *t-1*.

*** significant at 1% level; **significant at 5% level; *significant at 10% level

Figure 1. Investment rate over period by treatment

A) SPD



B) DPD

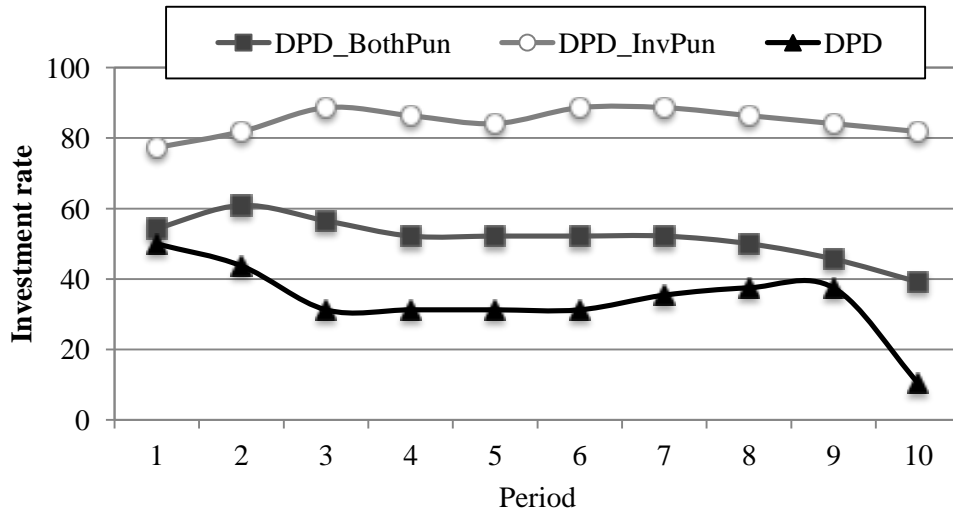
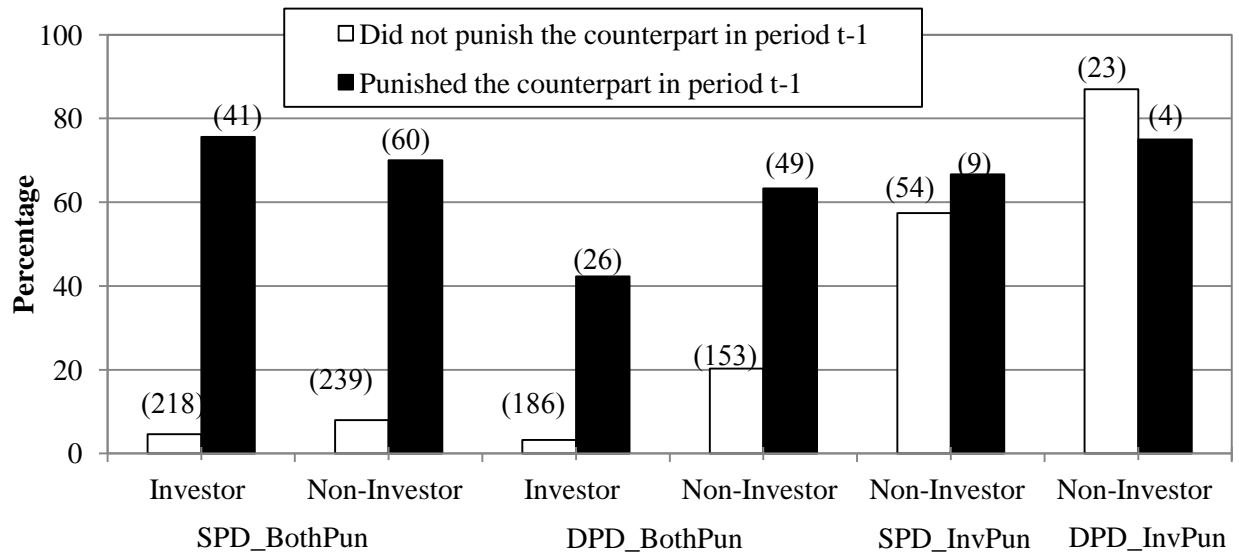


Figure 2. Frequency of receiving punishment in period t under each condition



*For SPD_InvPun and DPD_InvPun treatments, we count only the cases where subjects did not invest and the counterpart invested.

Numbers in the parenthesis are the numbers of observations in the corresponding cell.

Appendix A. Regression analysis of individual punishment decisions

To provide statistical evidence for the presence of retaliatory punishment, we define a variable ($PunExpen_{i,t}$): the amount subject i paid to punish the counterpart in the period t . This variable equals: (i) "0" if the subject did not punish the counterpart; (ii) "1" if the amount the subject paid to punish is positive and does not exceed four (punishment cost ratio of 1:3); and (iii) "2" if the amount paid to punish is greater than four. (We used three categories for punishment decisions in order to ensure that the number of observations in each cell is sufficiently large.) Using this as the dependent variable, we conducted a random individual effect ordered probit regression analysis of punishment decisions.

The independent variables include dummy variables for whether the subject i was punished in period $t-1$ ($Punished_{i,t-1}$ and $NotPunished_{i,t-1}$). We allow different coefficients for the case when the counterpart invested in period t and when the counterpart did not invest. The baseline is when subject i was not punished in period $t-1$ and the counterpart invested in period t . The regression results are reported in Table A-1. The coefficients of "... $Punished_{i,t-1}$ " are jointly significantly higher than the coefficients of "... $NotPunished_{i,t-1}$ " in each treatment (chi-square test, $p < 0.01$). This suggests that a subject is willing to pay more to punish her counterpart should the counterpart punish her in the previous period.¹⁶

To determine whether the occurrence of a loss plays a significant role in punishment decisions, we conducted a random individual effect ordered probit regression analysis of punishment decisions. The dependent variable is $PunExpen_{i,t}$. In the $SPD_BothPun$ treatment, the independent variables include a constant as well as three dummies for decision outcomes consisting of one's own decision and one's counterpart's decision: (i) oneself invested and one's counterpart did not invest ($Inv, NotInv$) $_{i,t}$; (ii) oneself did not invest and one's counterpart invested ($NotInv, Inv$) $_{i,t}$; (iii) neither invested ($NotInv, NotInv$) $_{i,t}$. The baseline is the case where both invested. In addition, for each decision outcome, we differentiate the case when the negative event happened ($NE_{i,t}=1$) and when the negative event did not happen ($NE_{i,t}=0$). The regression results are reported in Table A-2 Regression (1). We find the coefficients of ($Inv, NotInv$) $_{i,t} * NE_{i,t}$, ($NotInv, Inv$) $_{i,t} * NE_{i,t}$, and ($NotInv, NotInv$) $_{i,t} * NE_{i,t}$ are neither jointly nor individually significant (chi-square test, $p > 0.20$). This suggests that regardless of the investment decisions made by the subject or his counterpart, the occurrence of the negative event does not affect the subject's punishment decision.

¹⁶ Regression results from Table A-1 also show that the coefficient of " $CtpInv_{i,t} * Punished_{i,t-1}$ " is not significantly different from that of " $CtpNotInv_{i,t} * Punished_{i,t-1}$ " in both the $SPD_BothPun$ and $DPD_BothPun$ treatments (1.62 vs. 1.65, 1.97 vs. 1.99, chi-square test, $p > 0.90$). This suggests that when one is punished in the previous period, she is equally likely to punish the counterpart when the counterpart invested as when the counterpart did not invest in the current period. In contrast, the coefficient of " $CtpNotInv_{i,t} * NotPunished_{i,t-1}$ " is significantly positive in both the $SPD_BothPun$ (0.63) and $DPD_BothPun$ (1.55) treatments ($p < 0.05$), implying that when one is not punished in the previous period, she is more likely to punish non-investors than investors in the current period since the baseline case is $CtpInv_{i,t} * NotPunished_{i,t-1}$

In the *SPD_InvPun* treatment, since only the investors can punish only the non-investor counterpart, we include only the cases where the individual invested but her counterpart did not invest in the current period and $NE_{i,t}$ is the only independent variable. This regression result is reported in Table A-2 Regression (2). As shown in Regression (2), the coefficient of $NE_{i,t}$ is also not significant ($p>0.20$). Thus the regression results again suggest that the occurrence or absence of a negative event does not have a significant effect on one's punishment decision.

TableA-1. Random individual effect ordered probit regression analysis of punishment decisions in the *SPD_BothPun* and *DPD_BothPun* treatments

	Dependent variable: <i>PunExpen</i> _{<i>i,t</i>} (=0 if subject <i>i</i> did no punish the counterpart; =1 if subject <i>i</i> paid no more than 4 to punish the counterpart; =2 if subject <i>i</i> paid more than 4 to punish the counterpart)			
	<i>SPD_BothPun</i>		<i>DPD_BothPun</i>	
	Coef	s.e.	coef.	s.e.
<i>CtpInv</i> _{<i>i,t</i>} * <i>Punished</i> _{<i>i, t-1</i>}	1.62***	0.35	1.97***	0.38
<i>CtpNotInv</i> _{<i>i,t</i>} * <i>NotPunished</i> _{<i>i, t-1</i>}	0.63**	0.31	1.55***	0.30
<i>CtpNotInv</i> _{<i>i,t</i>} * <i>Punished</i> _{<i>i, t-1</i>}	1.65***	0.34	1.99***	0.34
<i>cut1_cons</i>	2.46	0.34	2.50	0.33
<i>cut2_cons</i>	3.42	0.37	3.31	0.37

Note: *** significant at 1% level; **significant at 5% level; *significant at 10% level

*CtpInv*_{*i,t*}=1 if *i*'s counterpart invested in period *t*; =0, o.w. ; *CtpNotInv*_{*i,t*}=1 if *i*'s counterpart did not invest in period *t*; =0, o.w. ; *NotPunished*_{*i, t-1*}=1 if *i* was not punished in period *t-1*; =0, o.w.;

*Punished*_{*i, t-1*}=1 if *i* was punished in period *t-1*; =0, o.w.

The baseline is when *i* was not punished in the period *t-1* and the counterpart invested in period *t*.

Table A-2. Random individual effect ordered probit regression analysis of punishment decisions in the SPD_BothPun and SPD_InvPun treatments

Dependent variable: <i>PunExpen</i> _{<i>i,t</i>} (=0 if subject <i>i</i> did no punish the counterpart; =1 if subject <i>i</i> paid no more than 4 to punish the counterpart; =2 if subject <i>i</i> paid more than 4 to punish the counterpart)		
	SPD_BothPun Regression (1)	SPD_InvPun* Regression (2)
	Coef. (s.e.)	Coef. (s.e.)
(Inv, NotInv) _{<i>i,t</i>}	1.07*** (0.30)	
(Inv, NotInv) _{<i>i,t</i>} *NE _{<i>i,t</i>}	-0.56 (0.48)	
(NotInv, Inv) _{<i>i,t</i>}	0.66** (0.30)	
(NotInv, Inv) _{<i>i,t</i>} *NE _{<i>i,t</i>}	0.05 (0.40)	
(NotInv, NotInv) _{<i>i,t</i>}	0.24 (0.32)	
(NotInv, NotInv) _{<i>i,t</i>} *NE _{<i>i,t</i>}	0.16 (0.30)	
NE _{<i>i,t</i>}		0.17 (0.30)
cut1_cons	2.19 (0.35)	-0.14 (0.22)
cut2_cons	3.10 (0.37)	0.89 (0.24)

Note: *** significant at 1% level; **significant at 5% level; *significant at 10% level
parenthesis is denoted as (*i*'s decision, *i*'s counterpart's decision)_{*i,t*} in period *t*;
NE_{*i,t*}=1 if the negative event happened in period *t*; =0, o.w.

*includes only the cases where the individual invested but his/her counterpart did not invest in the current period