

Goeschl, Timo; Jarke, Johannes

**Working Paper**

## Second vs. Third Party Punishment under Costly Monitoring: A New Experimental Method and Evidence

WiSo-HH Working Paper Series, No. 6

**Provided in Cooperation with:**

University of Hamburg, Faculty of Business, Economics and Social Sciences, WISO Research Lab

*Suggested Citation:* Goeschl, Timo; Jarke, Johannes (2013) : Second vs. Third Party Punishment under Costly Monitoring: A New Experimental Method and Evidence, WiSo-HH Working Paper Series, No. 6, Universität Hamburg, Fakultät für Wirtschafts- und Sozialwissenschaften, WiSo-Forschungslabor, Hamburg

This Version is available at:

<https://hdl.handle.net/10419/260412>

**Standard-Nutzungsbedingungen:**

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

**Terms of use:**

*Documents in EconStor may be saved and copied for your personal and scholarly purposes.*

*You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.*

*If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.*



Universität Hamburg  
DER FORSCHUNG | DER LEHRE | DER BILDUNG

FAKULTÄT  
FÜR WIRTSCHAFTS- UND  
SOZIALWISSENSCHAFTEN

# Second vs. Third Party Punishment under Costly Monitoring: A New Experimental Method and Evidence

Timo Goeschl  
Johannes Jarke

WiSo-HH Working Paper Series  
Working Paper No. 06



WiSo-HH Working Paper Series  
Working Paper No. 06

## **Second vs. Third Party Punishment under Costly Monitoring: A New Experimental Method and Evidence**

Timo Goeschl, Universität Heidelberg  
Johannes Jarke, Universität Hamburg

ISSN 2196-8128

Font used: „TheSans UHH“ / LucasFonts

Die Working Paper Series bieten Forscherinnen und Forschern, die an Projekten in Federführung oder mit der Beteiligung der Fakultät für Wirtschafts- und Sozialwissenschaften der Universität Hamburg tätig sind, die Möglichkeit zur digitalen Publikation ihrer Forschungsergebnisse. Die Reihe erscheint in unregelmäßiger Reihenfolge.

Jede Nummer erscheint in digitaler Version unter  
<https://www.wiso.uni-hamburg.de/de/forschung/working-paper-series/>

### **Kontakt:**

WiSo-Forschungslabor  
Von-Melle-Park 5  
20146 Hamburg

E-Mail: [experiments@wiso.uni-hamburg.de](mailto:experiments@wiso.uni-hamburg.de)

Web: <http://www.wiso.uni-hamburg.de/forschung/forschungslabor/home/>



# Second vs. Third Party Punishment under Costly Monitoring: A New Experimental Method and Evidence

Timo Goeschl\*

Johannes Jarke<sup>†</sup>

## Abstract

We experimentally evaluate a set of hypotheses regarding the effect of monitoring costs on the relative performance of second and third party punishment in the context of a social dilemma. By means of a novel method that mimicks information structures by constraints on the strategy space we reconcile the powerful «strategy method» with active information acquisition about the history of play. In line with recent behavioral and neuroscientific evidence we find second and third party punishment to be more similar than different. Both are less discriminate and hence generate weaker incentives for cooperation when monitoring is costly compared to a setting in which monitoring is free. Third parties do not exhibit more «careful» punishment behavior than second parties under costly monitoring, which sheds doubt on the «impartial spectator hypothesis». However, we also find subtle differences: The presence of monitoring costs leads subjects to withhold sanctioning more often as second parties than as third parties, and to punish indiscriminately more often as third parties than as second parties. We believe that the results further our understanding whether there is a common motivational structure underpinning both types of punishment or not. (*JEL* C92, C72, D03, D80)

## 1 Introduction

The resolution of social dilemmas can be greatly facilitated by the presence of persons that are willing to incur private costs in order to punish non-cooperative behavior (see e.g. Henrich et al., 2004; Gintis et al., 2005; Bowles & Gintis, 2011). Such persons fall in one of two classes: parties that are directly affected by the behavior (so-called «second parties») and parties that are materially unaffected (so-called «third parties»). However, sanctioning behavior by «vested» second parties and «by-standing» third parties might differ in critical ways. Understanding commonalities and differences between

---

\* Alfred-Weber-Institute of Economics, Universität Heidelberg. Bergheimer Strasse 20, D-69115 Heidelberg, Germany. Phone: +49 6221 54 8010. E-mail: goeschl@eco.uni-heidelberg.de.

<sup>†</sup> Corresponding author. School of Business, Economics and Social Sciences, Department of Socioeconomics, Universität Hamburg. Welckerstrasse 8, D-20354 Hamburg, Germany. Phone: +49 40 42838 8768. E-mail: johannes.jarke@wiso.uni-hamburg.de.

second and third party punishment is important because in many real social dilemmas only one of them is available.<sup>1</sup>

In a seminal study, Fehr & Fischbacher (2004b, FF in what follows) pioneered an experimental paradigm that combines a one-shot prisoners' dilemma with a costly punishment stage in order to compare costly sanctioning behavior by second parties and third parties in the context of a social dilemma. Each subject played the game in two conditions, a second party punishment condition (2PP) and a third party punishment condition (3PP), with randomly matched and unknown coplayers. In the first stage, all subjects could decide to cooperate with their matched partner or to defect. Subsequently, in the 2PP both players in each group could choose to punish the other. In the 3PP, each player had the opportunity to punish a player from another group. Punishment was costly to both the punishing subject and the punished subject, yet three times as costly for the latter compared to the former. In both conditions, subjects reported their punishment decisions using the *strategy method* (SM). Specifically, instead of directly responding to the first-stage actions with a single punishment action (the *direct-response method*, DRM), subjects specified a complete punishment strategy, i.e. a punishment level for each possible first-stage action profile, without being informed about the actual profile.<sup>2</sup> Apparently, the SM is substantially more powerful compared to the DRM with respect to the depth of data generated by an experiment because it not only elicits the factual path of play but also the counter-factual ones.<sup>3</sup>

FF found that both second and third parties engage in costly punishment, that is, spend some portion of their income to reduce others' income in spite of no prospect of a pecuniary return.<sup>4</sup> Furthermore, they find commonalities and differences between behavior by second and third parties, respectively. First, both second and third parties impose, on average, stronger punishments on defectors than on cooperators. Second, while punishments imposed on cooperators were approximately equal between second and third parties, the former imposed significantly stiffer punishments on defectors than the latter. This suggests that second parties could have a stronger preference to target sanctions on defectors than third parties.

This hypothesis has immediate implications for the realistic case in which a person that is in the position to punish is not perfectly informed, and hence uncertain about whether the target player cooperated or defected. In such situations the person can either decide to mete out punishments under

---

<sup>1</sup>Many legal institutions, for example, explicitly ban second party sanctions and replace them by sanctions determined by third parties. Another example are environmentally harmful behaviors, where there is no directly affected party at all.

<sup>2</sup>In the 2PP players announced a contingent punishment plan, indicating a punishment level for each of the target player's possible first-stage decisions, without being informed about the actual decision. Likewise, in the 3PP each third party, while being informed about the first-stage decision of the other member in their own group, had to indicate a punishment plan over four possible contingencies, one for each possible action profile in another group.

<sup>3</sup>However, the two methods may also produce different behaviors in the same game. See Brandts & Charness (2011) for a systematic comparison. Nevertheless, the evidence suggests that the differences are quite predictable.

<sup>4</sup>This result has been replicated in numerous studies. For overviews on costly second party punishment see Gächter & Herrmann (2009), Chaudhuri (2011) or Balliet et al. (2011). Henrich et al. (2004), Fehr & Fischbacher (2004a) and Nikiforakis & Mitchell (2014) summarize studies on third party punishment. We are aware of two further studies that directly compare second and third party punishment (Carpenter & Matthews, 2012; Leibbrandt & López-Pérez, 2012) but both rely on somewhat different designs compared to FF. We consider them in section 4.

uncertainty, bearing the risk of error,<sup>5</sup> or remedy the lack of information before deciding to punish (or not), possibly at a cost.<sup>6</sup> We call such remediation «monitoring». If second and third parties differ in their preference to target punishments on defectors, then we should observe differences in monitoring behavior, because the value of information about the target players behavior is different under the premise. Such differences have obvious and important ramifications for the relative efficacy of second and third party punishment in social dilemmas where monitoring is costly.<sup>7</sup> In the present paper we investigate this hypothesis (against competing hypotheses) within FF's paradigm and thereby develop a methodical innovation that allows for a simple and flexible implementation of monitoring into the SM. In doing so, we contribute to (i) the small but growing literature on the question whether and how second and third party punishment differ (see section 4), and (ii) the body of methods designed to elicit *strategies* in experimental games.

The implementation of monitoring within FF's paradigm poses a challenge because the crux of the SM is to keep subjects uninformed about (the prior history of) actual play. Informing subjects about actions chosen by the other players, or giving them the opportunity to acquire such information by themselves, renders the SM obsolete. We introduce a design feature that allows to implement monitoring in an *implicit* way, i.e. without actually informing subjects. The key idea behind this seemingly paradoxical feature is to mimic information structures by constraints on the strategy space: In FF the space of punishment strategies is (almost) unrestricted such that players eligible to punish can act *as if* they were perfectly informed about the target player's first-stage action. By imposing the same constraints on the strategy space that a more imperfect information structure would do, any degree of information imperfectness can be mimicked within the SM. We will expose this idea precisely in section 2.

We apply the method in an experiment that is built very closely on FF's. Specifically, in our experiment subjects in the role of a second or a third party had the (costly) option to condition their punishment on the behavior of the other player, i.e. treat cooperating and defecting behavior differently, or not. This is equivalent to a binary choice between the acquisition of perfect information on the target player's behavior («perfect monitoring»), and no information at all («no monitoring»). On the basis of a crossover design, we compare second and third party sanctioning behavior (within-subjects) at three different levels of monitoring costs (between-subjects), including zero. The zero monitor-

---

<sup>5</sup>There are two types of error, relative to her own preferences, a person  $P$  might perform under uncertainty: First,  $P$  could punish a person that she would not like to punish under certainty. Second,  $P$  could let a person escape unpunished that she would like to punish under certainty.

<sup>6</sup>Information on the target individual's actions is typically not immutably imperfect, but players are in a position to augment available information through costly effort. In other words, a player's choice whether to impose social sanctions on another player is frequently preceded by a *choice* on whether to invest resources into monitoring his or her actions. Such monitoring effort to overcome imperfect information on coplayers' actions is well known in a variety of economically relevant contexts, such as shared resource management (Ostrom, 1990; Rustagi et al., 2010), production teams and cooperatives (Acheson, 1975; Palmer, 1991; Kandel & Lazear, 1992; Dong & Dow, 1993; Craig & Pencavel, 1995), labor relations (Shapiro & Stiglitz, 1984; Kanemoto & MacLeod, 1991; Lazear, 1993), finance (Williamson, 1987; Armendáriz & Morduch, 2005) or neighborhood watch (Sampson et al., 1997), to name just a few.

<sup>7</sup>The implication is that monitoring costs have a more adverse affect on the disciplining function of third party punishment compared to second party punishment. See section 4.

ing cost condition is essentially a replication of FF. In the main conditions subjects therefore had not only a costly choice with respect to the intensity of punishment, but also a costly choice whether they wanted to «monitor».<sup>8</sup>

The experiment is designed to evaluate a set of hypotheses about commonalities and differences between second and third party behavior in settings where monitoring is costly. The «price effect hypothesis» applies to second and third parties alike and predicts that punishment will be less discriminate and hence generates weaker incentives for cooperation when monitoring is costly. The idea behind the hypothesis is that monitoring costs raise the price of discriminate punishment strategies, i.e. those that specify a different punishment level for a defector than for a cooperator, relative to the other strategies. Our second set of hypotheses consists of three alternative predictions regarding *differences* between second and third party behavior with respect to the effect predicted by the «price effect hypothesis». Specifically, the «vested interest hypothesis» predicts that the price effect is weaker among second parties than among third parties, the «impartial spectator hypothesis» predicts the opposite, and the «similarity hypothesis» predicts no difference.

Our findings broadly emphasize the commonalities between second and third party punishment. The «price effect hypothesis» is clearly supported and the «similarity hypothesis» defeats the two alternatives. At the same time, we find subtle differences that are not excluded by the «similarity hypothesis»: The presence of monitoring costs leads subjects to withhold sanctioning more often as second parties than as third parties, and to punish indiscriminately more often as third parties than as second parties.

In the remainder of the paper we proceed as follows. In section 2 we develop a simple theoretical framework to prove our claim that information structures can be mimicked by constraints on the strategy space. On that basis we describe the design and procedures of the experiment in section 3. In section 4 we derive the hypotheses with respect to the outcome of the experiment. We present the results in section 5 and conclude in section 6.

## 2 Theoretical framework

Consider the basic setup of the 2PP and 3PP. There are three players called  $A$  (the «first party»),  $B$  (the «second party») and  $C$  (the «third party»). We take the outcome of the first stage as given and model the second stage as a decision problem based on the standard state-space model of information (Hinikka, 1962; Aumann, 1976). Let  $\omega \equiv (\omega_A, \omega_B)$  be the action profile from the PD played by  $A$  and  $B$  in the first-stage, where  $\omega_i \in \{c, d\}$ ,  $c$  stands for «cooperate» and  $d$  for «defect». We consider the second stage from the perspective of person  $j$  that is eligible to punish person  $A$ , that is,  $j = B$  in the 2PP and  $j = C$  in the 3PP. The four possible first-stage action profiles constitute four possible worlds in which  $j$  can find herself. Formally, this is the state space  $\Omega \equiv \{(d, d), (d, c), (c, d), (c, c)\}$ .

An *information correspondence* for  $j$  is a map  $F_j : \Omega \rightarrow 2^\Omega \setminus \emptyset$  that associates to each state  $\omega \in \Omega$  a non-empty subset  $F_j(\omega) \subseteq \Omega$ , called  *$j$ 's information set* in state  $\omega$ .  $F_j(\omega)$  represents  $j$ 's informedness in that state: each state  $\omega \in F_j(\omega)$  is considered possibly true by  $j$ , and each state  $\omega \in \Omega \setminus F_j(\omega)$  is

---

<sup>8</sup>We explain later in the paper that monitoring costs are *not* the same as punishment costs.

certainly false. We impose the standard assumption that  $F_j$  is partitional, that is, each information function induces a partition  $\mathcal{F}_j$  of the state space  $\Omega$ .<sup>9</sup>

Let  $P_j$  be the set of punishment levels  $j$  is eligible to impose on  $A$ . Player  $j$ 's task is to choose a function  $p_j : \Omega \rightarrow P_j$ , called a *punishment strategy*, for which  $p_j(\omega) = p_j(\omega')$  whenever  $\omega \in F_j(\omega) \wedge \omega' \in F_j(\omega)$ . That is,  $j$  needs to specify a punishment level for each possible state under the constraint that she cannot choose different punishment levels in states that are indistinguishable to her. Thus,  $j$ 's degree of uninformedness imposes *restrictions* on the set of feasible punishment strategies.

If we denote the true state  $\bar{\omega}$ , the advantage of the strategy method (SM) over the direct-response method (DRM) can be precisely stated within this simple framework: The DRM elicits only  $p_j(\bar{\omega})$ , whereas the SM elicits the complete strategy  $p_j(\omega)$ . However, in all applications of the SM we are aware of there are minimal or no restrictions on the strategy space. Specifically, FF allow the specification of a different punishment level for each possible first-stage action performed by  $A$ : In the 3PP,  $p_C(\omega)$  is specified for each  $\omega \in \Omega$  and all four punishment levels are allowed to be different. In the 2PP  $p_B(\omega)$  is specified with the restrictions  $p_B(d, d) = p_B(d, c)$  and  $p_B(c, d) = p_B(c, c)$ , that is,  $B$  chooses a punishment level for each of  $A$ 's two possible first-stage actions.<sup>10</sup>

We now demonstrate that this application of the SM is equivalent to a setting in which  $A$ 's first-stage action can be *perfectly monitored* by  $j$  (i.e. a setting in which  $j$  receives a perfect signal on  $A$ 's first-stage action before imposing punishments). To see this, define  $S : \Omega \rightarrow 2^\Omega \setminus \emptyset$  that associates to each state  $\omega \in \Omega$  a subset  $S(\omega) \subseteq \Omega$  with  $\omega \in S(\omega)$ , called *signal*. Thus, any state produces a signal  $S(\omega)$  to  $j$  that contains the true state  $\omega$  and possibly other states. All states in  $S(\omega)$  are possible while all states in  $\Omega \setminus S(\omega)$  can be ruled out.<sup>11</sup> A signal is called *perfect* if it is a singleton that contains only the true state,  $S(\omega) = \{\omega\}$ . Otherwise it is called *imperfect* with the boundary case  $S(\omega) = \Omega$  that can be interpreted as «no signal». Because  $F_j$  is partitional and  $S(\omega)$  contains the true state,  $j$  deduces that the true state must be in both  $F_j(\omega)$  and  $S(\omega)$ . Thus, for all  $\omega \in S$  we have  $\hat{F}_j(\omega) = F_j(\omega) \cap S(\omega)$ , where  $\hat{F}_j$  is  $j$ 's updated information correspondence.<sup>12</sup>

Now, assume that  $j$  receives a perfect signal at the beginning of the second stage, before she metes out punishments. Then  $\hat{F}_j(\omega) = \{\omega\}$  for all  $\omega \in \Omega$ , or equivalently

$$\hat{\mathcal{F}}_j = \{\{(d, d)\}, \{(d, c)\}, \{(c, d)\}, \{(c, c)\}\}$$

<sup>9</sup>This is equivalent to the assumptions (i) that  $\omega \in F_j(\omega)$  holds for every  $\omega \in \Omega$  ( $j$  cannot be convinced of things that are false) and (ii) that  $\omega' \in F_j(\omega) \Rightarrow F_j(\omega') = F_j(\omega)$  ( $j$  uses the consistency or inconsistency of states with his information to make inferences about the state). See for example Osborne & Rubinstein (1994).

<sup>10</sup>Here FF are not entirely rigorous in applying the SM. Theoretically, a fully specified strategy would allow  $B$  also to condition her punishments on her own first-stage action. However, because we want to replicate FF's study as a baseline as close as possible, we follow their application.

<sup>11</sup>Put differently, any subset  $E \subseteq \Omega$  is called an *event*. Event  $E$  occurs in state  $\omega$  if and only if  $\omega \in E$ . Conversely, if  $E$  occurs the true state must be in  $E$ . The signal  $S(\omega)$  is an event that occurs in  $\omega$ .

<sup>12</sup>Thus, a signal in state  $\omega$  is *uninformative* if  $F_j(\omega) \subseteq S(\omega)$  because then  $\hat{F}_j(\omega) = F_j(\omega)$ , otherwise it is *informative*. In other words, a signal is informative if it allows  $j$  to rule out at least one state that she previously considered possible. Note that  $F_j(\omega)$  and  $S(\omega)$  can never be disjoint because they both contain the true state. Thus, a perfect signal is only uninformative if  $j$  already knew the true state. Because  $\hat{F}_j(\omega) \subseteq F_j(\omega)$  holds for all  $\omega \in \Omega$ , the induced information partition  $\hat{\mathcal{F}}_j$  is never *coarser* (less and bigger information sets) than  $\mathcal{F}_j$ , representing the fact that  $j$  cannot be less informed after receiving  $S(\omega)$ .



But this implies that there is no pair of states that is in the same information set, such that the set of feasible punishment strategies is not restricted by the information structure. Specifically, in the 3PP player  $C$  is free to choose a different punishment level  $p_C(\omega)$  in each of the four states. In the 2PP player  $B$  is free to specify a different punishment level for each of  $A$ 's two possible first-stage actions. This demonstrates that FF's implementation of the SM is *equivalent* to a naturally played (or DRM-played) game with perfect monitoring of  $A$ 's (and  $B$ 's) first-stage actions.

We can also show that this equivalence of information structures and restrictions of the strategy space naturally extends to cases with imperfect signals, including the case of no signal as defined above. Specifically, the imposition of the restriction  $p_j(\omega) = p_j(\omega')$  for some  $\omega \in \Omega$  and  $\omega' \in \Omega$  effectively mimics a particular information structure in which  $\omega$  and  $\omega'$  belong to the same information set. In the present paper, we are interested in the two extremes, perfect ignorance and perfect information about player  $A$ 's first-stage action. The latter is characterized formally above. Perfect ignorance in the 2PP is implemented by the restrictions  $p_B(c, a_B) = p_B(d, a_B)$  for any  $a_B \in \{c, d\}$ . Those restrictions on the set of punishment strategies are identical to the ones imposed by an information partition

$$\hat{\mathcal{F}}_B = \{\{(d, d), (c, d)\}, \{(d, c), (c, c)\}\}$$

i.e.  $B$  knows her own but not  $A$ 's first-stage action. Together with the restrictions  $p_B(a_A, d) = p_B(a_A, c)$  for any  $a_A \in \{c, d\}$  (see above) this implies that  $p_B(\omega)$  must be identical for all  $\omega \in \Omega$ . Likewise, in the 3PP the restrictions  $p_C(c, a_B) = p_C(d, a_B)$  for any  $a_B \in \{c, d\}$  on the set of punishment strategies are identical to the ones imposed by an information partition

$$\hat{\mathcal{F}}_C = \{\{(d, d), (c, d)\}, \{(d, c), (c, c)\}\}$$

i.e.  $C$  knows  $B$ 's but not  $A$ 's first-stage action.

Putting all this together, then, a decision to choose between a punishment strategy with or without the specified restrictions is equivalent to a decision between receiving a perfect signal and receiving no signal at all about  $A$ 's first-stage action. In other words, we can mimic a monitoring decision within the SM (and thus without actually informing  $j$  about  $A$ 's first-stage action) by allowing subjects to choose between a punishment strategy in which punishments are conditioned on  $a_A$  and a strategy in which punishments are *not* conditioned on  $a_A$ . We term the former strategy *discriminate* and the latter *indiscriminate*. Choosing a discriminate punishment strategy is equivalent to monitoring the target player before the application of punishments.

### 3 Experiment

Our basic experimental setup is essentially identical to FF. Each subject played a game consisting of a one-shot PD stage and a costly punishment stage in two conditions, a second party punishment condition (2PP) and a third party punishment condition (3PP), with randomly matched and unknown coplayers in random sequence.<sup>13</sup> In the first stage, all subjects could decide to cooperate by transfer-

<sup>13</sup>Randomization of the sequence was implemented by counter-balancing, i.e. in some sessions the 2PP was played before the 3PP, in others the 3PP before the 2PP. The assignment of subjects to those sequences was random.

ring their endowment of 10 tokens, in which case the experimenter tripled it, or to defect by keeping their endowment for themselves.

At the beginning of the second stage, each player received an additional endowment of 40 tokens in both conditions. Subsequently, in the 2PP both players in each group could choose to punish the other. In the 3PP, the first (second) player in a group had the opportunity to punish a player from another (a third) group.<sup>14</sup> A linear sanctioning technology was used:<sup>15</sup> any token spent by the punishing player subtracted three tokens from the punished player's payoff. Expenditures were restricted to integers and a maximum of 20 tokens.

Subjects reported their punishment decisions in the same way as in FF. Specifically, in the 2PP players announced a contingent punishment strategy, indicating a punishment level for each of the target player's possible transfer decisions, without being informed about the actual decision. Likewise, in the 3PP each third party, while being informed about the first-stage decision of the other member in their own group, had to indicate a punishment strategy over four possible contingencies.

Based on this setup we introduce the novel element of (costly) monitoring of the target player's first-stage action as described in the previous section. If a subject conditioned punishment on the target player's action, a payment of an exogenously set monitoring fee *in addition* to the applicable punishment costs fell due. The alternative was a punishment strategy in which the subject did not condition the punishment on the target player's action. This plan did not involve a fee over and above the cost of punishment.

We varied the level of the monitoring fee between-subjects at three levels: A fee of zero is an obvious starting point—it provides a direct replication of the PD experiment by FF and therefore a meaningful baseline treatment. In what follows, this treatment is labeled *M0*. Given the payoff structure of the game, we expected a fee of ten tokens close to prohibitive such that we have chosen this level as a reasonable boundary of the considered fee interval (treatment *M10*). In order to account for possible non-linearities we implemented another condition at the midpoint of the fee interval at five tokens (treatment *M5*). We should emphasize at this point that monitoring costs are *not* the same as punishment costs. On the basis of the theoretical framework developed in the previous section this is easy to see: monitoring costs are incurred *independent* from the realization of the state while the level of punishment costs depends on the state.

Participants were recruited from the general undergraduate student population using the online recruitment system ORSEE (Greiner, 2004). In total 134 subjects participated,<sup>16</sup> and under consideration of experimental efficiency (e.g. McClelland, 1997; List et al., 2011) we allocated approximately half of them to the midpoint cell *M5* (66 subjects), a third to maximum cell *M10* (46 subjects), and a sixth to the baseline cell *M0* (22 subjects), respectively. Of course, no subject participated in more than one session and monitoring fee treatment.

The experiment was conducted with a personal computer network and the software z-Tree (Fischbacher, 2007) at the experimental laboratory of the Alfred-Weber-Institute (AWI-Lab) at the Uni-

---

<sup>14</sup>This matching protocol rules out «tacit collusion» between third parties (see Fehr & Fischbacher, 2004b).

<sup>15</sup>See Casari (2005) on potential biases due to non-linear sanctioning technologies.

<sup>16</sup>Among the participants 49.3 percent were female and a share of 64.2 percent had never participated in a laboratory experiment before. The mean age was 21.8 years.

versity of Heidelberg. Upon entering the laboratory, subjects were randomly assigned to the computer terminals. Direct communication among them was not allowed for the duration of the entire session. Booths separated the participants visually, ensuring that they made their decisions anonymously and independently. Furthermore, subjects did not receive any information on the personal identity of any other participant, neither before nor while nor after the experiment.

At the beginning of the experiment, that is, before any decisions were made, subjects received detailed written instructions that explained the exact structure of the game and the procedural rules. The experiment was framed in a sterile way using neutral language and avoiding value laden terms in the instructions (see supplementary material). An advantage of our design is that all subjects in all conditions were confronted with *exactly* the same task and instructions that differed *only* in one number: the fee level. Participants had to answer a set of control questions individually at their respective seats in order to ensure comprehension of the rules. We did not start the experiment before all subjects had answered all questions correctly.

After that the exact timing of events was as follows. First, the subjects were randomly matched into groups of two. Then each subject made her or his decisions in the 2PP or 3PP, depending on the sequence.<sup>17</sup> After being informed about the payoffs, the experimenter announced that a second experiment will be conducted and distributed additional instructions that explained the differences to the first game. After being randomly re-matched into new groups, each subject made her or his decisions in the 3PP or 2PP, depending on the sequence. After being informed about the payoffs, subjects were asked to answer a short questionnaire while the experimenter prepared the payoffs. Subjects were then individually called to the experimenter booth, payed out (according to a random number matched to their decisions; no personal identities were used throughout the whole experiment) and dismissed. The whole experiment lasted approximately 90 minutes and subjects earned an average €18.46 (€0.10 per token earned), including a fixed show-up fee. Earnings exceed the local average hourly wage of a typical student job and can hence be considered meaningful to the participants.

## 4 Hypotheses

### 4.1 Baseline

Based on FF and a substantial body of further evidence on costly punishment (see the references in the introduction) we expect that the majority of second and third parties engages in costly punishment. Comparing second and third party punishment, FF's key results were as follows: First, both second and third parties intensely punished defectors whereas punishment of cooperators was present but infrequent and weak (but not absent). Second, punishment of defectors by second parties was stronger than third party punishment.<sup>18</sup> We expect to replicate those results in the baseline (*M0*) of our exper-

---

<sup>17</sup>Until the end of the first game, subjects did not know that another game will follow. This ambiguity was a deliberate design choice aimed at minimizing confounding effects of the second on the first game.

<sup>18</sup>For the readers convenience we repeat the underlying statistics to those results (in direct comparison to ours) in table 1 in parentheses. A further result was both cooperators and defectors punished defectors, but cooperators punished more frequently and imposed stronger sanctions. This result, however, was statistically not as solid as the others due to the small

iment. Furthermore, we expect that whenever a subject engages in costly punishment, (s)he chooses a discriminate punishment strategy, i.e. specifies a different punishment level for a defector than for a cooperator. This is because under costless monitoring discriminate and indiscriminate punishment strategies have the same price (zero) such that an arbitrarily weak preference to target punishment suffices to opt for the former.

## 4.2 Costly monitoring

With respect to the effect of monitoring costs we operationalize our hypotheses in terms of distributions over three «punishment types». The first type is one that decides not to punish at all, i.e.  $p_j(\omega) = 0$  for all  $\omega \in \Omega$ . This type is referred to as type  $N$ . Among the punishing subjects ( $p_j(\omega) > 0$  for some  $\omega \in \Omega$ ), there are two types. «Blind punishers» punish defectors and cooperators to the same extent, that is, they choose an indiscriminate strategy (as defined above). This type is denoted as  $B$ . The other type  $T$  are «targeted punishers» that choose a discriminate strategy (as defined above as well). The  $B$ -types and  $T$ -types together form the class of «punishers», denoted  $P$ .

Our first prediction, the «price effect hypothesis», applies to second and third parties alike: we expect less discriminate punishment when monitoring is costly. That is, we predict that the frequency of  $T$ -types will be lower in the costly monitoring conditions compared to the baseline condition in both the 2PP and the 3PP. As a direct consequence, the incentive to cooperate, defined as the difference between average punishment of defection and average punishment of cooperation, should be lower in the costly monitoring conditions compared to the baseline condition.

We base this hypothesis on two observations. First, previous evidence suggests that most people, both as a second and as a third party, target defectors when such targeting is costless (see references above). This suggests the existence of a preference for discriminate strategies which specify stronger punishments on defectors than on cooperators. Second, punishment behavior has been shown to be elastic with respect to incentives: increasing the cost of punishment reduces the supply of sanctions (Suleiman, 1996; Oosterbeek et al., 2004; Anderson & Putterman, 2006; Carpenter, 2007; Egas & Riedl, 2008). We conclude that since the price of discriminate punishment strategies relative to the other punishment strategies is higher in the costly monitoring conditions than in the baseline, they should be chosen less frequently in the former condition compared to the latter.

Conditional on the prediction with respect to the  $T$ -types the «price effect hypothesis» is consistent with any distribution over the  $N$ - and  $B$ -types: The flip-side of the prediction is that the *sum* of  $N$ - and  $B$ -types should be higher in the costly monitoring conditions compared to the baseline condition, but without additional assumptions we cannot say more. However, there are considerations that can sharpen our expectations. Consider a player with a preference to target punishments on defectors first. Under costless monitoring this player will behave as type  $T$ . Now, if a given monitoring fee renders the price of a discriminate punishment strategy prohibitive to this player, then refraining from punishment altogether (type  $N$ ) appears to be the intuitively natural response. However, there is also evidence on pure «money burning» or «nastiness» (Zizzo & Oswald, 2001; Zizzo, 2003; Abbink & Sadrieh, 2009;

---

number of defectors. Furthermore, it is not of primary interest for our purposes here. The interested reader finds this result replicated in appendix A.

Abbink & Herrmann, 2011), that is, costly «punishment» without any history of previous interaction on which it could be conditioned upon. This implies that we also could expect to observe at least some *B*-types under costly monitoring.

Our second set of alternative hypotheses is about *differences* between second and third party behavior with respect to the effect predicted by the «price effect hypothesis». Specifically, FF found that cooperators were punished, on average, equally strong by second and third parties, respectively, but defectors were punished markedly stronger by second parties compared to third parties.<sup>19</sup> This suggests that second parties might have a stronger preference for discriminate punishment strategies than third parties, which would imply that the average second party has a higher reservation price for discriminate strategies than the average third party.<sup>20</sup> A motivational underpinning could be that direct reciprocity (e.g. Gouldner, 1960; Goranson & Berkowitz, 1966; Rabin, 1993; Fehr & Gächter, 1998; Dufwenberg & Kirchsteiger, 2004; Falk & Fischbacher, 2006) can be a concern for second parties but not for third parties, because there is no history of previous interaction between a third party and the target player. Hence, third parties might be just less interested in the target player's first-stage behavior. Under this premise of different reservation prices we should observe that a given raise of the monitoring fee «crowds out» relatively more third parties than second parties from the class of *T*-types. As a consequence, a given raise of the monitoring fee should have a stronger (negative) effect on the incentive to cooperate (as defined above) in the 3PP than in the 2PP.

We term this prediction «vested interest hypothesis» and place it against two alternative hypotheses. First, precisely because third parties cannot be motivated by direct reciprocity («revenge»), they are often suspected to be less egocentric, more impartial, and motivated more strongly by morality and ethical norms than second parties (Fehr et al., 2002; Fehr & Fischbacher, 2004b,a). This idea is exemplified by Adam Smith's (1759) notion of an «impartial spectator» and is embodied in many institutions that rely on third parties to determine sanctions (e.g. courts and juries).<sup>21</sup> If this hypothesis is true then we should observe the opposite of the prediction by the «vested interest hypothesis»: a given raise of the monitoring fee should «crowd out» relatively more second parties than third parties from the class of *T*-types. As a consequence, a given raise of the monitoring fee should have a stronger (negative) effect on the incentive to cooperate in the 2PP than in the 3PP. We term this the «impartial spectator hypothesis».

Our final alternative hypothesis holds that punishment behavior by second and third parties is actually sufficiently similar such that we should observe no significant differences between them when adjusting to monitoring costs. We term it «similarity hypothesis». There are good arguments in favor of the «similarity hypothesis». Theoretically, a number of possible motivations to punish at a cost have been proposed in the literature, such as distributional preferences (e.g. Fehr & Schmidt,

---

<sup>19</sup>In FF's experiment, second parties spent on average 8.4 tokens to punish defectors and 0.7 tokens to punish cooperators. Third parties spent between 1.4 and 3.0 tokens to punish defectors and between 0.6 and 0.7 tokens to punish cooperators.

<sup>20</sup>Put differently, the average second party has a stronger aversion to place punishments under uncertainty (i.e. indiscriminately) because errors (punishing a cooperators; letting a defector escape unpunished) are experienced worse.

<sup>21</sup>However, there is evidence that third parties can have their own egocentric agendas and biases (Ross et al., 1977; Babcock et al., 1995) and evidence from third-party-driven sanctioning institutions (e.g. Thatcher, 1891; Kennedy, 1997) shows that third parties sometimes diverge strongly from the «impartial spectator» ideal.

1999; Bolton & Ockenfels, 2000; Charness & Rabin, 2002), spite (e.g. Levine, 1998), or norms (e.g. López-Pérez, 2008). Except direct reciprocity, all of them can in principle motivate both second and third parties alike.<sup>22</sup> Empirically, recent studies comparing second and third party punishment using different experimental paradigms from FF's and ours appear to support the notion that they are more similar than different. Leibbrandt & López-Pérez (2012) use a sequence of ten (dictator) distribution games, each followed by a punishment (by second parties in one condition and by third parties in the other) stage, respectively. Their key results are twofold. First, and in contrast to FF, third parties punish as intensely as second parties. Second, and contrary to folk wisdom (the «impartial spectator hypothesis»), third party punishment does not exhibit a more normative or impartial pattern than second party punishment.<sup>23</sup> Evaluating a set of preference theories by means of a classification analysis, the authors find that inequity aversion is the best explanation of both second and third party punishment. Interpreting their and previous results, Leibbrandt & López-Pérez (2012) conclude that second and third party punishment may actually be more similar than previously thought.<sup>24</sup>

Evidence from neuroscience lends further support to the «similarity hypothesis». Buckholtz et al. (2008) scanned subjects with functional magnetic resonance imaging (fMRI) while they determined the appropriate punishment for crimes that varied in perpetrator responsibility and crime severity. Their interpretation of the evidence grants activity in the right dorsolateral prefrontal cortex a key role in third party punishment, and since the same prefrontal region has previously been shown to be involved in second party punishment (de Quervain et al., 2004) and other two-player cooperation games (Sanfey et al., 2003; Singer et al., 2004, 2006; Delgado et al., 2005; King-Casas et al., 2005; Knoch et al., 2006; Spitzer et al., 2007), the authors suggest that the cognitive processes supporting second and third party punishment may be supported by a common neural mechanism.<sup>25</sup>

---

<sup>22</sup>In repeated games, prospects of pecuniary returns in the future can be a motive for punishment by either party, albeit somewhat more remotely for third parties (Kreps et al., 1982; Fudenberg & Maskin, 1986; Kandori, 1992). However, this motive is ruled out by the one-shot design of FF's and our experiment.

<sup>23</sup>Specifically, third parties tend to punish those dictators that choose allocations in which they become the richest party, even if that choice is joint payoff maximizing or Pareto efficient or equitable in the target pair. Furthermore, third parties even punish players who cannot make a choice at all, especially if the latter are richer than the third party. See also López-Pérez & Leibbrandt (2011). However, in the context of a public good game Carpenter & Matthews (2012) find that third parties intervene primarily to promote efficiency. One possible explanation for the mixed evidence is differences in experimental designs. The strong support for distributional motives in López-Pérez & Leibbrandt (2011) and Leibbrandt & López-Pérez (2012) might have to do with the use of distribution games which render distributional aspects strongly salient, whereas a public good game or a PD may trigger different motivations. Similar effects have been observed in the context of second party punishment and other economic games were different experimental games highlight different motivations (see e.g. Falk et al., 2003, 2005, 2008; Engelmann & Strobel, 2004; Fowler et al., 2005; Dawes et al., 2007; Johnson et al., 2009). In fact, in their PD experiment Fehr & Fischbacher (2004b) observe a significant fraction of third parties punishing mutual defection, which is incompatible with distributional theories, given the payoff structure in their experiment (the third party was already the subject with the highest payoff after mutual defection in the target group).

<sup>24</sup>There is one further study doing the comparison in a repeated games setting: Carpenter & Matthews (2009) study a public good game with punishment played repeatedly by the same group of four players. In one treatment only «in-group punishment» (second party punishment) is possible, in another there is also the opportunity to punish players in a different group (third party punishment). The added complexity arising from repetition renders it difficult to compare it to the other studies and to draw inferences with respect to motivations.

<sup>25</sup>See Seymour et al. (2007) and Fehr (2009) for reviews of the neuroscientific literature on punitive behavior.

**Table 1:** Intensity and frequency of second and third party punishment in the baseline condition.

Target player is	TP punishment		
	SP punishment	Target player's coplayer is	
		Cooperator	Defector
Defector	7.6 (8.4) <i>72.7% (66.7%)</i>	6.0 (3.1) <i>72.7% (58.7%)</i>	2.3 (1.4) <i>45.5% (32.6%)</i>
Cooperator	0.6 (0.7) <i>13.6% (8.3%)</i>	1.4 (0.6) <i>22.7% (15.2%)</i>	1.1 (0.7) <i>22.7% (15.2%)</i>

Average (per subject) expenditure for punishment points in regular letters. Relative frequency of punishing subjects in italics. The results of Fehr and Fischbacher (2004) are in parantheses.

## 5 Results

We organize the presentation of the experimental results in the following way. We begin by presenting our results from the baseline treatment that replicates the study by FF. Against this benchmark, we then describe the core results of the experiment.

### 5.1 Baseline

The results of our baseline treatment are shown in table 1 (FF's results are shown in parentheses).<sup>26</sup> As evident from the table, our first two predictions are clearly supported: Almost three quarters of the subjects punished defectors (around one fifth punished cooperators), and defectors were punished significantly stronger on average than cooperators by both second and third parties.<sup>27</sup> The «targeting» prediction is supported as well: all punishing second parties (16 out of 16) and all but one punishing third party (15 out of 16) have chosen discriminate punishment strategies, respectively.

The prediction that second party punishment is stronger than third party punishment is, in line with Leibbrandt & López-Pérez (2012), supported only partially. Table 1 shows that second parties spent more to punish defectors on average, but the difference is statistically significant only in case the target player's coplayer also defected (Wilcoxon signed rank test,  $p = .0004$ ) and insignificant in case the target player's coplayer cooperated ( $p = .2171$ ). This is due to the fact that third party punishment, especially of unilateral defectors, is notably stronger in our experiment than in FF.

**Table 2:** Intensity and frequency of second and third party punishment in the costly monitoring conditions.

Target player is	TP punishment					
	SP punishment		Target player's coplayer is			
	<i>M5</i>	<i>M10</i>	Cooperator	<i>M5</i>	<i>M10</i>	Defector
Defector	3.3 <i>31.8%</i>	1.9 <i>26.1%</i>	2.6 <i>36.4%</i>	2.8 <i>41.3%</i>	1.6 <i>25.8%</i>	1.7 <i>30.4%</i>
Cooperator	1.1 <i>16.7%</i>	0.3 <i>13.0%</i>	1.3 <i>22.7%</i>	1.5 <i>32.6%</i>	1.0 <i>19.7%</i>	0.6 <i>17.4%</i>

Average (per subject) expenditure for punishment points in regular letters. Relative frequency of punishing subjects in italics.

## 5.2 Costly monitoring

Table 2 extends table 1 to the costly monitoring conditions. Our first prediction, the «price effect hypothesis», is supported by the data. Consider the aggregate results shown in table 2 first. On the one hand, defectors are still punished stronger than cooperators on average by both second and third parties in the costly monitoring conditions.<sup>28</sup> On the other hand, sanctions are weaker and less directed to defectors if monitoring is costly.<sup>29</sup> Hence, defecting got less costly as monitoring costs rise. To see this, observe that the incentive to cooperate (as defined above) decreases substantially as monitoring gets costly, both in the 2PP and the 3PP (see also table 4): In the 2PP the incentive decreases by 4.8 points (69.0 percent) from the baseline to *M5* and by 5.5 points (78.3 percent) from the baseline to *M10*. Both differences are statistically significant using Mann-Whitney rank-sum tests ( $p = .0000$  and  $p = .0000$ , respectively). The incentive to cooperate in the 3PP decreases by 3.3 points (71.6 percent,  $p = .0000$ ) from the baseline to *M5* and by 3.3 points (71.9 percent,  $p = .0000$ ) from the baseline to *M10* in case the target player's coplayer cooperated.<sup>30</sup>

Now consider the distributions of our «punishment types» shown at the margins of table 3, where

<sup>26</sup>We do not find any significant effects of order, that is, whether the 2PP is played before the 3PP or the other way around (see appendix B), such that we pool the data from both sequences in order to increase statistical power. But note that all key results from this section also hold in each sequence individually, as shown in appendix B.

<sup>27</sup>Wilcoxon signed-rank tests yield  $p = .0001$  in case of second party punishment,  $p = .0002$  in case of third party punishment when the target player's coplayer cooperated, and  $p = .0988$  (marginal significance) in case target player's coplayer also defected.

<sup>28</sup>In the *M5* treatment, Wilcoxon signed-rank tests yield  $p = .0063$  in case of second party punishment,  $p = .0009$  in case of third party punishment when the target player's coplayer cooperated, and  $p = .0143$  in case target player's coplayer cooperated. In the *M10* treatment, the test results are  $p = .0144$ ,  $p = .0054$ , and  $p = .0082$ , respectively.

<sup>29</sup>In both the 2PP and the 3PP punishment of cooperators is independent from the monitoring fee (Kruskal-Wallis tests,  $p = .7563$  in the 2PP,  $p = .7955$  in the 3PP if the target player's coplayer defects, and  $p = .5995$  if the coplayer cooperates as well). Punishment of defectors, however, is affected by the fee ( $p = .0002$  in the 2PP,  $p = .0051$  in the 3PP if the target player's coplayer cooperates, and  $p = .2909$  if the coplayer defects as well).

<sup>30</sup>In case the coplayer defected as well, third party punishment decreases by 0.6 points (49.3 percent,  $p = .1354$ ) from the baseline to *M5* and by close to zero points (2.4 percent,  $p = .4916$ ) from the baseline to *M10*.



**Table 3:** «Transition matrices» of behavioral types across the 2PP and the 3PP.

Condition <i>M0</i>		3PP Type			
2PP Type	<i>N</i>	<i>T</i>	<i>B</i>	Margin	
<i>N</i>	(4) 18.2%	(2) 9.1%	(0) 0.0%	(6) 27.3%	
<i>T</i>	(2) 9.1%	(13) 59.1%	(1) 4.5%	(16) 72.7%	
<i>B</i>	(0) 0.0%	(0) 0.0%	(0) 0.0%	(0) 0.0%	
Margin	(6) 27.3%	(15) 68.2%	(1) 4.5%	(22) 100%	

  

Condition <i>M5</i>		3PP Type			
2PP Type	<i>N</i>	<i>T</i>	<i>B</i>	Margin	
<i>N</i>	(37) 56.1%	(3) 4.5%	(4) 6.1%	(44) 66.7%	
<i>T</i>	(2) 3.0%	(7) 10.6%	(5) 7.6%	(14) 21.2%	
<i>B</i>	(2) 3.0%	(1) 1.5%	(5) 7.6%	(8) 12.1%	
Margin	(41) 62.1%	(11) 16.7%	(14) 21.2%	(66) 100%	

  

Condition <i>M10</i>		3PP Type			
2PP Type	<i>N</i>	<i>T</i>	<i>B</i>	Margin	
<i>N</i>	(25) 54.3%	(3) 6.5%	(6) 13.4%	(34) 73.9%	
<i>T</i>	(0) 0.0%	(3) 6.5%	(3) 6.5%	(6) 13.4%	
<i>B</i>	(1) 2.2%	(2) 4.3%	(3) 6.5%	(6) 13.4%	
Margin	(26) 56.5%	(8) 17.4%	(12) 26.1%	(46) 100%	

Absolute frequencies in parantheses.

**Table 4:** Aggregate relative performance of second and third party punishment at different monitoring fee levels.

Condition	Incentive 3PP				
	Incentive 2PP	Target player's coplayer is		2P overhead	
		Cooperator	Defector	Cooperator	Defector
<i>M0</i>	7.0	4.6	1.1	2.4	5.9
<i>M5</i>	2.2	1.3	0.6	0.9	1.6
<i>M10</i>	1.5	1.3	1.1	0.2	0.4

"Incentive" means punishment of cooperation subtracted from punishment from defection. "SP overhead" means "incentive" in the 3PP subtracted from "incentive" in the 2PP.

the column margin represents the 2PP and the row margin the 3PP (ignore the interior cells of the table for the moment). From those distributions it is readily apparent that variations in the cost of monitoring are associated with different patterns of punishment behavior, which are fully consistent with the «price effect hypothesis»: First, as shown by the *N*-rows and *N*-columns in table 3, there is a higher share of second and third parties that do not punish at all in the costly monitoring conditions compared to the baseline condition.<sup>31</sup>

Second, in addition to the result on the propensity to punish at all, our design explicitly allows to study punitive behavior of those individuals that opt to punish independently from the target player's first stage behavior. The presence of those types are shown in the *B*-rows and *B*-columns in table 3. Among both second and third parties there is a greater fraction of blind punishers in the costly monitoring conditions compared to the baseline.<sup>32</sup>

Third, as evident from the *T*-rows and *T*-columns in table 3, the share of targeted punishers is notably smaller in the costly monitoring conditions compared to the baseline in both second and third parties.<sup>33</sup>

We now turn to our second set of hypotheses. Since the frequency of *T*-types is not significantly different between second and third parties in the baseline condition, the «vested interest hypothesis» reduces to the prediction that the *T*-types should more frequent among the second parties than among

<sup>31</sup>The hypothesis that the frequency of *N*-types and monitoring costs are independent can be rejected (Kruskal-Wallis test,  $p = .0007$  in the 2PP,  $p = .0171$  in the 3PP). Pairwise, the differences between *M0* and *M5* (Mann-Whitney tests,  $p = .0013$  in the 2PP,  $p = .0048$  in the 3PP) and *M0* and *M10* ( $p = .0003$  in the 2PP,  $p = .0248$  in the 3PP) are significant, the differences between *M5* and *M10* are not ( $p = .4140$  in the 2PP,  $p = .5539$  in the 3PP).

<sup>32</sup>In the baseline the frequency of blind punishers is not significantly different from zero (Wilcoxon test,  $p = .3173$  in the 3PP; in the 2PP the frequency is strictly zero), whereas in the costly monitoring conditions they are ( $p = .0047$  in the 2PP-*M5*,  $p = .0002$  in the 3PP-*M5*,  $p = .0143$  in the 2PP-*M10*,  $p = .0005$  in the 3PP-*M10*). Pairwise, the differences between *M0* and *M5* (Mann-Whitney tests,  $p = .0886$  in the 2PP,  $p = .0734$  in the 3PP) and *M0* and *M10* ( $p = .0782$  in the 2PP,  $p = .0359$  in the 3PP) are marginally significant, the differences between *M5* and *M10* are not ( $p = .8851$  in the 2PP,  $p = .5495$  in the 3PP).

<sup>33</sup>In both cases the hypothesis that the frequency of *T*-types and monitoring costs are independent can be rejected (Kruskal-Wallis tests,  $p = .0001$ , respectively). Pairwise, the differences between *M0* and *M5* (Mann-Whitney tests,  $p = .0000$  in the 2PP,  $p = .0000$  in the 3PP) and *M0* and *M10* ( $p = .0000$  in the 2PP,  $p = .0000$  in the 3PP) are significant, the differences between *M5* and *M10* are not ( $p = .2689$  in the 2PP,  $p = .9203$  in the 3PP).

the third parties in the costly monitoring conditions. The «impartial spectator hypothesis» predicts the opposite, and the «similarity hypothesis» predicts that the frequency of *T*-types does not differ between second and third parties in the costly monitoring conditions.

The data shown in table 3 supports the «similarity hypothesis». Comparing the frequency of *T*-types between the 2PP and the 3PP in both costly monitoring conditions pooled together using a Wilcoxon test yields  $p = .8185$  (for the two conditions separately  $p = .3657$  in *M5* and  $p = .4795$  in *M10*). Thus, the decrease of *T*-types as monitoring gets costly is not significantly different between second and third parties.

This is mirrored by a comparison of the incentives to cooperate across the 2PP and 3PP, which are shown in table 4. For clarity, we subtract the incentive to cooperate in the 3PP from the one in the 2PP and call the resulting value «2P overhead».<sup>34</sup> If this overhead is zero, second and third party punishment are equally strong (in the sense of generating the same incentive to cooperate); the larger the overhead, the stronger is second party punishment relative to third party punishment (where a negative value indicates that third party punishment is stronger). The values are shown in the fourth and fifth column of table 4. They permit two observations: First, all overheads are positive, i.e. second party punishment always generates stronger incentives to cooperate. Second, the overheads are strictly decreasing in the monitoring fee level, i.e. the incentives to cooperate in the 2PP and 3PP, respectively, converge towards equality as the monitoring fee increases. However, most differences are not statistically significant, such that the «similarity hypothesis» cannot be rejected without doubt.<sup>35</sup>

Summing up so far, there are clearly a number of commonalities between second and third party punishment under costly monitoring. However, we close the section by stressing a subtle difference that is not excluded by the «similarity hypothesis». First, the increase of *N*-types as monitoring gets costly is somewhat stronger in second than in third parties: the frequency of *N*-types is apparently not different between second and third parties in the baseline condition, but higher in second parties than in third parties when monitoring is costly.<sup>36</sup>

Second, the increase of *B*-types as monitoring gets costly is somewhat stronger in third than in second parties. In the baseline condition, the frequency of *B*-types is not significantly different from zero in both second and third parties, but strictly positive and higher in third parties than in second parties when monitoring is costly.<sup>37</sup> Thus, as a response to monitoring costs, second parties have a

<sup>34</sup>Note that due to our experimental design harnessing the strategy method and a within-subject comparison of second and third party behavior, this «difference in differences» is identified at the *individual* level.

<sup>35</sup>Neither difference in the fifth column is significant (Mann-Whitney tests,  $p = .2343$  for *M0* vs. *M5*,  $p = .3548$  for *M0* vs. *M10*, and  $p = .7671$  for *M5* vs. *M10*). In the sixth column, *M0* vs. *M5* ( $p = .0001$ ) and *M0* vs. *M10* ( $p = .0000$ ) are significant, *M5* vs. *M10* is not ( $p = .2970$ ).

<sup>36</sup>Comparing the frequency of *N*-types between the 2PP and the 3PP in both costly monitoring conditions pooled together using a Wilcoxon signed-rank test yields  $p = .0164$  (for the two conditions separately  $p = .3657$  in *M5* and  $p = .0114$  in *M10*). As an alternative demonstration, a dummy indicating treatment (monitoring is costly) is significantly positively correlated with a dummy indicating type *N* (Kendall's  $\tau_b = 0.2995$ ,  $p = .0006$ ), whereas this correlation is somewhat stronger for second parties ( $\tau_b = 0.3245$ ,  $p = .0002$ ) than for third parties ( $\tau_b = 0.2421$ ,  $p = .0053$ ).

<sup>37</sup>Comparing the frequency of *B*-types between the 2PP and the 3PP in both costly monitoring conditions pooled together using a Wilcoxon test yields  $p = .0143$  (for the two conditions separately  $p = .0833$  in *M5* and *M10*, respectively). A dummy

slightly stronger tendency to switch to non-punishment while third parties have a slightly stronger tendency to switch to blind punishment.

Those aggregate differences between the 2PP and the 3PP are driven by individuals that behave differently when they find themselves in the role of a second or third party, respectively. This is shown by the interior cells of table 3: across all monitoring fee conditions, 27.6 percent (37 out of 134) of the subjects change their behavioral type between the roles. Those subjects are represented by the off-main-diagonal cells in table 3, where each cell illustrates a different transition path. The key insight from this data is that (i) there are disproportionately more subjects (19) that punish blindly (type *B*) as a third but not as a second party than the other way around (6 subjects), and that (ii) there are disproportionately more subjects (18) that do not punish at all (type *N*) as a second but not as a third party than the other way around (7 subjects).

## 6 Conclusion

In this paper we investigated how the relative performance of costly second and third party punishment in the context of a social dilemma is affected by the presence of monitoring costs. On the way, we suggested a simple and flexible trick to implement monitoring into the «strategy method», and applied it within the seminal experimental paradigm designed by Fehr & Fischbacher (2004b) to compare second and third party punishment. In line with recent behavioral (Leibbrandt & López-Pérez, 2012) and neuroscientific evidence (Buckholtz et al., 2008) our results emphasize the commonalities between second and third party punishment. However, we also find subtle differences.

We believe that the results further our understanding whether there is a common motivational structure underpinning both types of punishment or not. So far, the experimental evidence on this question has been less than conclusive. Our findings show that the responses of second and third party punishment to variations in the cost of monitoring yield very similar patterns. This is not to be interpreted as evidence that both types of punishment share an identical motivational structure. But a similar motivational structure would be consistent with this finding. Notably, the finding that third parties do not exhibit more «careful» punishment behavior than second parties under costly monitoring sheds doubt on the «impartial spectator hypothesis».

## References

- Abbink, K. & Herrmann, B. (2011). The moral cost of nastiness. *Economic Inquiry*, 49(2), 631–633.
- Abbink, K. & Sadrieh, A. (2009). The pleasure of being nasty. *Economics Letters*, 105(3), 306–308.
- Acheson, J. M. (1975). The lobster fiefs: Economic and ecological effects of territoriality in the maine lobster industry. *Human Ecology*, 3(3), 183–207.
- Anderson, C. M. & Putterman, L. (2006). Do non-strategic sanctions obey the law of demand? The demand for punishment in the voluntary contribution mechanism. *Games and Economic Behavior*, 54(1), 1–24.
- Armendáriz, B. & Morduch, J. (2005). *The Economics of Microfinance*. Cambridge, MA, USA: MIT Press.

---

indicating treatment (monitoring is costly) is significantly positively correlated with a dummy indicating type *B* (Kendall's  $\tau_b = 0.2066$ ,  $p = .0174$ ), whereas this correlation is somewhat weaker for second parties ( $\tau_b = 0.1514$ ,  $p = .0818$ ) than for third parties ( $\tau_b = 0.1724$ ,  $p = .0473$ ).

- Aumann, R. J. (1976). Agreeing to disagree. *Annals of Statistics*, 4(6), 1236–1239.
- Babcock, L., Loewenstein, G., Issacharoff, S., & Camerer, C. (1995). Biased judgments of fairness in bargaining. *American Economic Review*, 85(5), 1337–1343.
- Balliet, D., Mulder, L., & van Lange, P. A. M. (2011). Reward, punishment, and cooperation: A meta-analysis. *Psychological Bulletin*, 137(4), 594–615.
- Bolton, G. E. & Ockenfels, A. (2000). ERC: A theory of equity, reciprocity, and competition. *American Economic Review*, 90(1), 166–193.
- Bowles, S. & Gintis, H. (2011). *A Cooperative Species: Human Reciprocity and its Evolution*. Princeton, NJ, USA: Princeton University Press.
- Brandts, J. & Charness, G. (2011). The strategy versus the direct-response method: A first survey of experimental comparisons. *Experimental Economics*, 14(3), 375–398.
- Buckholtz, J. W., Asplund, C. L., Dux, P. E., Zald, D. H., Gore, J. C., Jones, O. D., & Marois, R. (2008). The neural correlates of third-party punishment. *Neuron*, 60(5), 930–940.
- Carpenter, J. (2007). The demand for punishment. *Journal of Economic Behavior & Organization*, 62(4), 522–542.
- Carpenter, J. P. & Matthews, P. H. (2009). What norms trigger punishment? *Experimental Economics*, 12(3), 272–288.
- Carpenter, J. P. & Matthews, P. H. (2012). Norm enforcement: Anger, indignation, or reciprocity? *Journal of the European Economic Association*, 10(3), 555–572.
- Casari, M. (2005). On the design of peer punishment experiments. *Experimental Economics*, 8(2), 107–115.
- Charness, G. & Rabin, M. (2002). Understanding social preferences with simple tests. *Quarterly Journal of Economics*, 117(3), 817–869.
- Chaudhuri, A. (2011). Sustaining cooperation in laboratory public goods experiments: A selective survey of the literature. *Experimental Economics*, 14(1), 47–83.
- Craig, B. & Pencavel, J. (1995). Participation and productivity: A comparison of worker cooperatives and conventional firms in the Plywood industry. *Brookings Papers on Economic Activity: Microeconomics*, 1995, 121–174.
- Dawes, C. T., Fowler, J. H., Johnson, T., McElreath, R., & Smirnov, O. (2007). Egalitarian motives in humans. *Nature*, 446, 794–796.
- de Quervain, D. J.-F., Fischbacher, U., Treyer, V., Schellhammer, M., Schnyder, U., Buck, A., & Fehr, E. (2004). The neural basis of altruistic punishment. *Science*, 305(5688), 1254–1258.
- Delgado, M. R., Frank, R. H., & Phelps, E. A. (2005). Perceptions of moral character modulate the neural systems of reward during the trust game. *Nature Neuroscience*, 8, 1611–1618.
- Dong, X.-Y. & Dow, G. K. (1993). Monitoring costs in Chinese agricultural teams. *Journal of Political Economy*, 101(3), 539–553.
- Dufwenberg, M. & Kirchsteiger, G. (2004). A theory of sequential reciprocity. *Games and Economic Behavior*, 47(2), 268–298.
- Egas, M. & Riedl, A. (2008). The economics of altruistic punishment and the maintenance of cooperation. *Philosophical Transactions of the Royal Society: Biological Sciences*, 275(1637), 871–878.
- Engelmann, D. & Strobel, M. (2004). Inequality aversion, efficiency, and maximin preferences in simple distribution experiments. *American Economic Review*, 94(4), 857–869.
- Falk, A., Fehr, E., & Fischbacher, U. (2003). On the nature of fair behavior. *Economic Inquiry*, 41(1), 20–26.
- Falk, A., Fehr, E., & Fischbacher, U. (2005). Driving forces behind informal sanctions. *Econometrica*, 73(6), 2017–2030.
- Falk, A., Fehr, E., & Fischbacher, U. (2008). Testing theories of fairness—intentions matter. *Games and Economic Behavior*, 62(1), 287–303.
- Falk, A. & Fischbacher, U. (2006). A theory reciprocity. *Games and Economic Behavior*, 54(2), 293–315.
- Fehr, E. (2009). Social preferences and the brain. In P. W. Glimcher, C. F. Camerer, E. Fehr, & R. A. Poldrack (Eds.), *Neuroeconomics: Decision Making and the Brain* (pp. 215–232). London: Academic Press.
- Fehr, E. & Fischbacher, U. (2004a). Social norms and human cooperation. *Trends in Cognitive Science*, 8(4), 185–190.
- Fehr, E. & Fischbacher, U. (2004b). Third-party punishment and social norms. *Evolution and Human Behavior*, 25(2), 63–87.

- Fehr, E., Fischbacher, U., & Gächter, S. (2002). Strong reciprocity, human cooperation, and the enforcement of social norms. *Human Nature*, 13(1), 1–25.
- Fehr, E. & Gächter, S. (1998). Reciprocity and economics: The economic implications of Homo Reciprocans. *European Economic Review*, 42(3-5), 845–859.
- Fehr, E. & Schmidt, K. M. (1999). A theory of fairness, competition, and cooperation. *Quarterly Journal of Economics*, 114(3), 817–868.
- Fischbacher, U. (2007). z-tree: Zurich toolbox for ready-made economic experiments. *Experimental Economics*, 10(2), 171–178.
- Fowler, J. H., Johnson, T., & Smirnov, O. (2005). Human behaviour: Egalitarian motive and altruistic punishment. *Nature*, 433, E1.
- Fudenberg, D. & Maskin, E. (1986). The folk theorem in repeated games with discounting or with incomplete information. *Econometrica*, 54(3), 533–554.
- Gächter, S. & Herrmann, B. (2009). Reciprocity, culture and human cooperation: Previous insights and a new cross-cultural experiment. *Philosophical Transactions of the Royal Society: Biological Sciences*, 364(1518), 791–806.
- Gintis, H., Bowles, S., Boyd, R., & Fehr, E., Eds. (2005). *Moral Sentiments and Material Interests. The Foundations of Cooperation in Economic Life*. Cambridge, USA: MIT Press.
- Goranson, R. E. & Berkowitz, L. (1966). Reciprocity and responsibility reactions to prior help. *Journal of Personality and Social Psychology*, 3(2), 227–232.
- Gouldner, A. W. (1960). The norm of reciprocity: A preliminary statement. *American Sociological Review*, 25(2), 161–178.
- Greiner, B. (2004). *The Online Recruitment System ORSEE 2.0—A Guide for the Organization of Experiments in Economics*. Working Paper Series in Economics 10, University of Cologne, Cologne, Germany.
- Henrich, J., Boyd, R., Bowles, S., Camerer, C., Fehr, E., & Gintis, H., Eds. (2004). *Foundations of Human Sociality: Economic Experiments and Ethnographic Evidence from Fifteen Small-Scale Societies*. Oxford, UK: Oxford University Press.
- Hinikka, J. (1962). *Knowledge and Belief: An Introduction to the Logic of the Two Notions*. Ithaca (NY, USA): Cornell University Press.
- Johnson, T., Dawes, C. T., Fowler, J. H., McElreath, R., & Smirnov, O. (2009). The role of egalitarian motives in altruistic punishment. *Economics Letters*, 102(3), 192–194.
- Kandel, E. & Lazear, E. P. (1992). Peer pressure and partnerships. *Journal of Political Economy*, 100(4), 801–817.
- Kandori, M. (1992). Social norms and community enforcement. *Review of Economic Studies*, 59(1), 63–80.
- Kanemoto, Y. & MacLeod, W. B. (1991). The theory of contracts and labor practices in Japan and the United States. *Managerial and Decision Economics*, 12(2), 159–170.
- Kennedy, R. (1997). *Race, Crime, and the Law*. New York, NY, USA: Pantheon.
- King-Casas, B., Tomlin, D., Anen, C., Camerer, C. F., Quartz, S. R., & Montague, P. R. (2005). Getting to know you: Reputation and trust in a two-person economic exchange. *Science*, 308(5718), 78–83.
- Knoch, D., Pascual-Leone, A., Meyer, K., Treyer, V., & Fehr, E. (2006). Diminishing reciprocal fairness by disrupting the right prefrontal cortex. *Science*, 314(5800), 829–832.
- Kreps, D. M., Milgrom, P., Roberts, J., & Wilson, R. (1982). Rational cooperation in the finitely repeated prisoners' dilemma. *Journal of Economic Theory*, 27(2), 245–252.
- Lazear, E. P. (1993). Labor economics and the psychology of organizations. *Journal of Economic Perspectives*, 5(2), 89–110.
- Leibbrandt, A. & López-Pérez, R. (2012). An exploration of third and second party punishment in ten simple games. *Journal of Economic Behavior & Organization*, 84(3), 753–766.
- Levine, D. K. (1998). Modeling altruism and spitefulness in experiments. *Review of Economic Dynamics*, 1(3), 593–622.
- List, J. A., Sadoff, S., & Wagner, M. (2011). So you want to run an experiment, now what? Some simple rules of thumb for optimal experimental design. *Experimental Economics*, 14(4), 439–457.
- López-Pérez, R. (2008). Aversion to norm-breaking: A model. *Games and Economic Behavior*, 64(1), 237–267.
- López-Pérez, R. & Leibbrandt, A. (2011). The dark side of altruistic third-party punishment. *Journal of Conflict Resolution*,

- 55(5), 761–784.
- McClelland, G. H. (1997). Optimal design in psychological research. *Psychological Methods*, 2(1), 3–19.
- Nikiforakis, N. & Mitchell, H. (2014). Mixing the carrots with the sticks: Third party punishment and reward. *Experimental Economics*, 17(1), 1–23.
- Oosterbeek, H., Sloof, R., & van de Kuilen, G. (2004). Cultural differences in ultimatum game experiments: Evidence from a meta-analysis. *Experimental Economics*, 7(2), 171–188.
- Osborne, M. J. & Rubinstein, A. (1994). *A course in game theory*. Cambridge, Massachusetts, USA: MIT Press.
- Ostrom, E. (1990). *Governing the commons: The evolution of institutions for collective action*. Cambridge, UK: Cambridge University Press.
- Palmer, C. T. (1991). Kin-selection, reciprocal altruism, and information sharing among Maine lobstermen. *Ethology and Sociobiology*, 12(3), 221–235.
- Rabin, M. (1993). Incorporating fairness into game theory and economics. *American Economic Review*, 83(3), 525–548.
- Ross, L., Greene, D., & House, P. (1977). The 'false consensus effect': An egocentric bias in social perception and attribution processes. *Journal of Experimental Social Psychology*, 13(3), 279–301.
- Rustagi, D., Engel, S., & Kosfeld, M. (2010). Conditional cooperation and costly monitoring explain success in forest commons management. *Science*, 330, 961–965.
- Sampson, R. J., Raudenbush, S. W., & Earls, F. (1997). Neighborhoods and violent crime: A multilevel study of collective efficacy. *Science*, 277(5328), 918–924.
- Sanfey, A. G., Rilling, J. K., Aronson, J. A., Nystrom, L. E., & Cohen, J. D. (2003). The neural basis of economic decision-making in the ultimatum game. *Science*, 300(5626), 1755–1758.
- Seymour, B., Singer, T., & Dolan, R. (2007). The neurobiology of punishment. *Nature Reviews Neuroscience*, 8, 300–311.
- Shapiro, C. & Stiglitz, J. E. (1984). Equilibrium unemployment as a worker discipline device. *American Economic Review*, 74(3), 433–444.
- Singer, T., Kiebel, S. J., Winston, J. S., Dolan, R. J., & Frith, C. D. (2004). Brain responses to the acquired moral status of faces. *Neuron*, 41(4), 653–662.
- Singer, T., Seymour, B., O'Doherty, J. P., Stephan, K. E., Dolan, R. J., & Frith, C. D. (2006). Empathic neural responses are modulated by the perceived fairness of others. *Nature*, 439, 466–469.
- Smith, A. (1759). *The Theory of Moral Sentiments*. London, England: A. Millar.
- Spitzer, M., Fischbacher, U., Herrnberger, B., Grön, G., & Fehr, E. (2007). The neural signature of norm compliance. *Neuron*, 56(1), 185–196.
- Suleiman, R. (1996). Expectations and fairness in a modified ultimatum game. *Journal of Economic Psychology*, 17(5), 531–554.
- Thatcher, C. A. (1891). The failure of the jury system. *North American Review*, 153(417), 246–249.
- Williamson, S. D. (1987). Costly monitoring, loan contracts, and equilibrium credit rationing. *Quarterly Journal of Economics*, 102(1), 135–146.
- Zizzo, D. J. (2003). Money burning and rank egalitarianism with random dictators. *Economics Letters*, 81(2), 263–266.
- Zizzo, D. J. & Oswald, A. J. (2001). Are people willing to pay to reduce others' incomes? *Annales D'Économie et de Statistique*, 63/64, 39–65.

**Table 5:** Who sanctions defectors in the baseline condition.

Punisher is	Target player is	SP punishment	TP punishment	
			Target player's coplayer is	
			Cooperator	Defector
Cooperator	Defector	8.3 (9.2) <i>77.8% (69.0%)</i>	6.0 (3.7) <i>81.3% (67.7%)</i>	2.7 (1.7) <i>56.3% (35.5%)</i>
Defector	Defector	4.3 (2.7) <i>50.0% (50.0%)</i>	6.0 (1.9) <i>50.0% (40.0%)</i>	1.2 (0.9) <i>16.7% (26.7%)</i>

Average (per subject) expenditure for punishment points in regular letters. Relative frequency of punishing subjects in italics. The results of Fehr and Fischbacher (2004) are in parantheses.

**Table 6:** Replication of table 1 separated by order of the 2PP and 3PP.

Target player is	SP punishment	TP punishment	
		Target player's coplayer is	
		Cooperator	Defector
Defector	6.9 / 8.9 <i>p = .7038</i>	6.4 / 5.3 <i>p = .3328</i>	2.1 / 2.5 <i>p = .8224</i>
Cooperator	0.9 / 0.0 <i>p = .1700</i>	1.7 / 0.8 <i>p = .4569</i>	1.4 / 0.8 <i>p = .4571</i>

Average (per subject) expenditure for punishment points in regular letters. On the left-hand side of the slash the 2PP is followed by the 3PP, on the right-hand side of the slash the 3PP is followed by the 2PP. The p-value of a Mann-Whitney test for order effects in italics.

## A Who sanctions defectors?

In footnote 18 we referred to a further result of Fehr & Fischbacher (2004b): both cooperators and defectors punished defectors, but cooperators punished more frequently and imposed stronger sanctions. As evident from table 5 (again, for the readers convenience we repeat the results of FF in parentheses) we also do find that cooperators punish defectors somewhat stronger than defectors in the baseline condition, although none of the three differences is statistically significant for the same reason as in the FF study (Mann-Whitney rank sum tests,  $p \geq .1823$ ).

## B Robustness to effects of order

Table 6 replicates table 1 showing average punishment levels separately for the 2PP-3PP sequence (left-hand side of the slash) and the 3PP-2PP sequence (right-hand side of the slash), respectively. Results of Mann-Whitney tests for order effects are shown in italics. Evidently, none of the differences are statistically significant. The main results of section 5.1 turn out to be robust to effects of order as well: First, second parties punished defectors significantly stronger than cooperators both in the 2PP-3PP sequence ( $p = .0036$ ) and the 3PP-2PP sequence ( $p = .0138$ ). Likewise, third parties punished



**Table 7:** Replication of table 2 separated by order of the 2PP and 3PP.

Target player is	TP punishment					
	SP punishment		Target player's coplayer is			
	<i>M5</i>	<i>M10</i>	Cooperator	Defector	<i>M5</i>	<i>M10</i>
Defector	4.7 / 1.6 <i>p = .0054</i>	1.5 / 2.3 <i>p = .5757</i>	3.6 / 1.4 <i>p = .1666</i>	2.7 / 3.1 <i>p = .9505</i>	2.2 / 0.8 <i>p = .0304</i>	1.7 / 1.7 <i>p = .8807</i>
Cooperator	1.5 / 0.6 <i>p = .0506</i>	0.3 / 0.4 <i>p = .7332</i>	2.0 / 0.5 <i>p = .0720</i>	1.3 / 1.8 <i>p = .6701</i>	1.7 / 0.2 <i>p = .0132</i>	0.8 / 0.3 <i>p = .2402</i>

Average (per subject) expenditure for punishment points in regular letters. On the left-hand side of the slash the 2PP is followed by the 3PP, on the right-hand side of the slash the 3PP is followed by the 2PP. The p-value of a Mann-Whitney test for order effects in italics.

unilateral defectors significantly stronger than cooperators both in the 2PP-3PP sequence ( $p = .0017$ ) and the 3PP-2PP sequence ( $p = .0487$ ).<sup>38</sup>

Second, punishment of defectors by second parties is not weaker than third party punishment. In the 3PP-2PP sequence second parties spent significantly more in punishment than third parties in all cases ( $p = .0487$  in case the target player's coplayer cooperated,  $p = .0190$  in case the target player's coplayer also defected), in the 2PP-3PP sequence only in case the target player's coplayer also defected ( $p = .0092$ ), whereas there is no significant difference in case the target player's coplayer cooperated ( $p = .8951$ ).

Table 7 replicates table 2, again showing the respective numbers separately for the 2PP-3PP sequence and the 3PP-2PP sequence, respectively, and with results of Mann-Whitney tests for order effects shown in italics. Punishment levels are also not significantly different between the two sequences in the *M10* condition. In the *M5* condition punishment levels are overall weaker in the 3PP-2PP than in the 2PP-3PP condition.<sup>39</sup> However, since our analysis focuses on *differences* (punishment of defectors vs. punishment of cooperators; punishment by second parties vs. punishment by third parties) none of the key results reported in section 5.2 are affected by this.

In particular, inspecting tables 8 and 9, which replicate table 3 separately for the 2PP-3PP sequence and the 3PP-2PP sequence, respectively, the key patterns of behavioral types described in section 5.2 are clearly present in both sequences.<sup>40</sup> First, there is a higher share of *N*-types in the costly monitoring

<sup>38</sup>As before, defectors whose coplayer defected as well are punished at most weakly stronger than cooperators: The difference is marginally significant in the 3PP-2PP sequence ( $p = .0848$ ) and insignificant in the 2PP-3PP sequence ( $p = .3992$ ).

<sup>39</sup>The difference is mainly due to one session with very low punishment levels. Since we do not observe such low levels in any other session, and no differences between sequences in the *M10* condition, we believe that this is a non-systematic outlier.

<sup>40</sup>We omit hypothesis tests because the number of observations in some cells are too small to meaningfully apply them. What we want to show here is that the key qualitative results reported in section 5.2 are apparent in both sequences. For readers interested in more detail we refer the data-set and the program code of the analysis that is provided in the supplementary material.

**Table 8:** Replication of table 3 with the data from the 2PP-3PP sequence only.

Condition <i>M0</i>		3PP Type		
2PP Type	<i>N</i>	<i>T</i>	<i>B</i>	Margin
<i>N</i>	(3) 21.4%	(2) 14.3%	(0) 0.0%	(5) 35.7%
<i>T</i>	(0) 0.0%	(9) 64.3%	(0) 0.0%	(9) 64.3%
<i>B</i>	(0) 0.0%	(0) 0.0%	(0) 0.0%	(0) 0.0%
Margin	(3) 21.4%	(11) 78.6%	(0) 0.0%	(14) 100%

  

Condition <i>M5</i>		3PP Type		
2PP Type	<i>N</i>	<i>T</i>	<i>B</i>	Margin
<i>N</i>	(16) 44.4%	(0) 0.0%	(2) 5.6%	(18) 50.0%
<i>T</i>	(2) 5.6%	(5) 13.9%	(5) 13.9%	(12) 33.3%
<i>B</i>	(2) 5.6%	(0) 0.0%	(4) 11.1%	(6) 16.7%
Margin	(20) 55.6%	(5) 13.9%	(11) 30.6%	(36) 100%

  

Condition <i>M10</i>		3PP Type		
2PP Type	<i>N</i>	<i>T</i>	<i>B</i>	Margin
<i>N</i>	(15) 57.7%	(2) 7.7%	(3) 11.5%	(20) 76.9%
<i>T</i>	(0) 0.0%	(1) 3.8%	(2) 7.7%	(3) 11.5%
<i>B</i>	(0) 0.0%	(1) 3.8%	(2) 7.6%	(3) 11.5%
Margin	(15) 57.7%	(4) 15.4%	(7) 26.9%	(26) 100%

Absolute frequencies in parantheses.

**Table 9:** Replication of table 3 with the data from the 3PP-2PP sequence only.

Condition <i>M0</i>		3PP Type		
2PP Type	<i>N</i>	<i>T</i>	<i>B</i>	Margin
<i>N</i>	(1) 12.5%	(0) 0.0%	(0) 0.0%	(1) 12.5%
<i>T</i>	(2) 25.0%	(4) 50.0%	(1) 12.5%	(7) 87.5%
<i>B</i>	(0) 0.0%	(0) 0.0%	(0) 0.0%	(0) 0.0%
Margin	(3) 37.5%	(4) 50.0%	(1) 12.5%	(8) 100%

  

Condition <i>M5</i>		3PP Type		
2PP Type	<i>N</i>	<i>T</i>	<i>B</i>	Margin
<i>N</i>	(21) 70.0%	(3) 10.0%	(2) 6.7%	(26) 86.7%
<i>T</i>	(0) 0.0%	(2) 6.7%	(0) 0.0%	(2) 6.7%
<i>B</i>	(0) 0.0%	(1) 3.3%	(1) 3.3%	(2) 6.7%
Margin	(21) 70.0%	(6) 20.0%	(3) 10.0%	(30) 100%

  

Condition <i>M10</i>		3PP Type		
2PP Type	<i>N</i>	<i>T</i>	<i>B</i>	Margin
<i>N</i>	(10) 50.0%	(1) 5.0%	(3) 15.0%	(14) 70.0%
<i>T</i>	(0) 0.0%	(2) 10.0%	(1) 5.0%	(3) 15.0%
<i>B</i>	(1) 5.0%	(1) 5.0%	(1) 5.0%	(3) 15.0%
Margin	(11) 55.0%	(4) 20.0%	(5) 25.0%	(20) 100%

Absolute frequencies in parantheses.

conditions compared to the baseline condition. Second, among both second and third parties there is a greater fraction of *B*-types in the costly monitoring conditions compared to the baseline. Third, the share of *T*-types is notably smaller in the costly monitoring conditions compared to the baseline in both second and third parties.

With respect to the differences between second and third party behavior, the two key results reported in the second part of section 5.2 are also apparent in both sequences individually. First, there are disproportionately more subjects (12 in the 2PP-3PP sequence, 7 in the 3PP-2PP sequence) that punish blindly (type *B*) as a third but not as a second party than the other way around (3 in both the 2PP-3PP and the 3PP-2PP sequence, respectively). Second, there are disproportionately more subjects (9 in both the 2PP-3PP and the 3PP-2PP sequence, respectively) that do not punish at all (type *N*) as a second but not as a third party than the other way around (4 in the 2PP-3PP sequence, 3 in the 3PP-2PP sequence). As a result, the increase of *N*-types as monitoring gets costly is somewhat stronger in second than in third parties, whereas the increase of *B*-types as monitoring gets costly is somewhat stronger in third than in second parties.