

International Journal of the Commons
Vol. 8, no 1 February 2014, pp. 47–78
Publisher: Igitur publishing
URL: <http://www.thecommonsjournal.org>
URN: NBN:NL:UI:10-1-116081
Copyright: content is licensed under a Creative Commons Attribution 3.0 License
ISSN: 1875-0281

Implementing punishment and reward in the public goods game: the effect of individual and collective decision rules

Nynke van Miltenburg

Department of Sociology/ICS, Utrecht University, The Netherlands
n.vanmiltenburg@uu.nl

Vincent Buskens

Department of Sociology/ICS, Utrecht University and Erasmus School of Law, Erasmus
University Rotterdam, The Netherlands
v.buskens@uu.nl

Davide Barrera

Department of Culture, Politics, and Society, University of Turin, Italy, and Department of
Sociology/ICS, Utrecht University, The Netherlands
davide.barrera@unito.it

Werner Raub

Department of Sociology/ICS, Utrecht University, The Netherlands
w.raub@uu.nl

Abstract: Punishments and rewards are effective means for establishing cooperation in social dilemmas. We compare a setting where actors *individually* decide whom to sanction with a setting where sanctions are only implemented when actors *collectively* agree that a certain actor should be sanctioned. Collective sanctioning decisions are problematic due to the difficulty of reaching consensus. However, when a decision is made collectively, perverse sanctioning (e.g. punishing high contributors) by individual actors is ruled out. Therefore, collective sanctioning decisions are likely to be in the interest of the whole group. We employ a laboratory experiment where subjects play Public Goods Games with opportunities for punishment or reward that is implemented either by an individual, a majority, or unanimously. For both punishment and reward, contribution levels are higher in the individual than the majority condition, and

higher under majority than unanimity. Often, majority agreement or unanimity was not reached on punishments or rewards.

Keywords: Collective decision rule, conditional cooperation, public goods game, reward

Acknowledgements: Comments of and discussions with colleagues of our Utrecht group “Cooperation in Social and Economic Relations”, and suggestions of anonymous reviewers are gratefully acknowledged. The paper is part of the project “The Feasibility of Cooperation under Various Sanctioning Institutions,” funded by the Netherlands Organization for Scientific Research (NWO, grant 400-09-159). Earlier drafts were presented in 2011 at the Symposium of the International Sociological Association (ISA) Research Committee 45 on Rational Choice, at the Fourth Conference of the European Network of Analytical Sociologists (ENAS), at the Nuffield Workshop on Experimental Research on Social Dilemmas, and in 2012 at the Conference on Design and Dynamics of Collective Action.

1. Introduction

A public good is characterized by non-excludability: once it is produced, all actors can enjoy its benefits regardless of their contribution to the provision of the good (Olson 1965 [1971]). Since public good provision is costly, this implies a tension between the individual and collective interest. While mutual cooperation leads to the best possible group outcome, individuals have an incentive to free-ride on the contributions of others.

Contributions to public goods can be supported by positive or negative peer sanctions, that is, the opportunity for actors to reward or punish each other. Experimental research established that high contributions are maintained when sanctioning is possible (Yamagishi 1986; Ostrom et al. 1992; Fehr and Gächter 2000, 2002; Sefton et al. 2007; Balliet et al. 2011). However, a challenge for research and policy is to design institutions that best enable heterogeneous actors to enforce cooperation (Ostrom 2010, 2012). In this respect, which method of implementing sanctions is most successful in increasing contributions remains an open question (Gächter and Thöni 2011). For example, it is unclear whether contributions are higher when the decision of whom to sanction is made individually or when it is made collectively. A sanctioning system with an individual decision rule (IDR) is a system in which every actor individually decides whom to sanction and pays the associated costs. A sanctioning system with a collective decision rule (CDR) is a system in which sanctions are executed only when multiple actors agree and pay the cost of sanctioning.

In real-life public good problems, actors often employ a sanctioning institution with a CDR. For example, Ostrom (1990) and Veszteg and Narhetali (2010) describe small communities where group members successfully enforce collective action through collective sanctioning decisions. Typically, members

of the community regularly meet to identify free-riders and decide upon their punishment, for example in a vote. Also, in international cooperation, nations use collective sanctioning decision rules to ensure provision of global public goods such as international security and economic stability. Sanctioning decisions are usually taken by a variant of majority voting. Unanimity voting is uncommon, because it gives every individual nation the opportunity to veto a sanction, thereby making collective organizations ineffective decision makers (www.europa.eu, www.un.org).

So far, there is limited experimental research comparing the effect of sanctioning through IDRs and CDRs in public good problems.¹ Casari and Luini (2009) find that, compared to an IDR, contributions to public goods are higher when punishment is only carried out if at least two out of four actors punish the fifth member of their group. Thus, they consider only one CDR, and do not compare positive and negative sanctions. This leaves a number of unresolved issues, in which the current paper provides further insight.

First, it is unclear how the effect of a CDR on contribution depends upon the proportion of actors required to agree for a sanction to be implemented. On the one hand, the higher the proportion required, the less likely it will be that a sufficient number of actors agrees on the necessity of sanctioning and is willing to incur the associated costs (cf. Buchanan and Tullock 1962). Thus, while under an IDR all desired sanctions are carried out by definition, under CDRs there is a higher chance that free-riders remain unpunished or contributors unrewarded. On the other hand, under an IDR individuals might decide to use sanctions in ways that hurt contribution and thereby result in decreasing payoffs for the group, i.e. to reward free-riders or to punish contributors (Casari and Luini 2009; Ellingsen et al. 2012). Consequently, the more actors collectively agree that a certain group member should be sanctioned, the higher the chance that this sanction will be in the collective interest, that is, in accordance with enforcing contributions to the public good. In the current paper, we address the effect of the required proportion of consenting actors on contribution levels by comparing contributions under an IDR to a CDR for which majority and a CDR for which unanimity is required.

Second, theoretical arguments and empirical results on punishment cannot be straightforwardly generalized to reward. For example, to maintain cooperation rewards have to be repeatedly allocated to contributors. Conversely, the mere threat of punishment can be enough to deter free-riding (Dari-Mattiacci and De

¹ Numerous experimental studies employ other forms of collective peer sanctioning decisions, but do not compare IDRs with CDRs. Decker et al. (2003) examine the effect of implementing a subset of punishment proposals while actors share punishment costs. In Ertan et al. (2009), Sutter et al. (2010), and Botelho et al. (2005) subjects vote on whether to allow individual peer sanctioning. Andreoni and Gee (2012); Guillen et al. (2006); Kamei et al. (2011); Markussen et al. (2011); Putterman et al. (2011); Traulsen et al. (2012) and Tyran and Feld (2006) let subjects collectively decide on implementing an institution which automatically punishes free-riders.

Geest 2010). This suggests that punishments and rewards may differ in efficiency. Empirically, it has been shown that punishments and rewards might differ also in terms of efficacy (e.g. Sefton et al. 2007; Choi and Ahn 2013). We therefore study decision rules for assigning both punishment and reward.

The effects of the decision rules on macro-behavior such as aggregate contribution levels depend on assumptions about the micro-motives of individual actors (cf. Gächter and Thöni 2011). For example, these effects depend on which proportion of actors is willing to sanction, who is likely to be targeted, and how sanctions influence contribution decisions. We summarize existing knowledge on individual behavior in the PGG with sanctions. Subsequently, we apply this to predict macro-level behavior in the PGG with different decision rules, and with punishment or reward. We thus assess through which mechanisms our empirical extensions could result in different contribution levels between sanctioning systems.

The paper is structured as follows. In the theory section, we review the literature on behavior in public good problems with opportunities for sanctioning. Subsequently, we develop hypotheses about contribution and sanctioning behavior, and on how this behavior of individuals translates in different contribution levels under IDRs and CDRs. Individual-level and macro-level hypotheses are tested in an experiment where individual, majority, and unanimity decision rules for punishing and rewarding are employed in an incentivized manner.

2. Theory

2.1 The Public Goods Game

The linear Public Goods Game (PGG; also called Voluntary Contribution Mechanism, e.g. Isaac and Walker 1988) is used as a model of public good problems. It is played by n actors. All actors i receive an endowment w . They simultaneously and independently decide whether to keep this endowment for themselves or contribute an amount $g_i \in [0, w]$ to a “group account”. The total amount contributed by all n actors together, $g = \sum g_i$, is multiplied by a number m , with $1 < m < n$, and mg is divided equally among all actors. Because $m < n$, the individual return obtained from the amount contributed to the group account is smaller than when it would have been kept to oneself ($mg_i/n < g_i$). Therefore, when the PGG is played once under standard game-theoretic assumptions – that is, when actors are rational in maximizing utility and selfish in that utility equals own payoff – contributing nothing is a dominant strategy, yielding the highest utility regardless what others do. This results in the unique Nash equilibrium of no contributions. However, since $m > 1$ the joint group outcome $nw - g + mg$ is maximized when everybody contributes the full endowment. Every player would then be better off compared to when all contribute nothing ($mw > w$). Thus, individually rational behavior leads to a Pareto-suboptimal outcome, making the PGG a social dilemma (Dawes 1980).

2.2 Behavior in the PGG

The prediction of complete free-riding is typically refuted in experimental research employing the PGG. Instead, contributions averaging 50% of the endowment are consistently observed in one-shot PGGs (Walker and Halloran 2004; Kocher et al. 2007). Also in repeated PGGs where group composition changes after each round, as in our experiment, subjects initially contribute around 50% on average. However, in subsequent rounds contributions gradually decline to very low levels (Ledyard 1995).

Research explaining this declining contribution pattern focuses on non-standard utility as an alternative behavioral assumption. It has been empirically established that actors in the PGG can be classified in two main preference types (Ostrom 2000; Fehr and Gintis 2007; Ones and Putterman 2007). Actors of the first type are rational and selfish *free-riders* who never contribute to the public good. Actors of the second type are *conditional cooperators* who contribute more, the more they expect others to contribute (see Gächter 2007; Chaudhuri 2011 for an overview of empirical evidence). These actors are assumed to derive utility from reciprocating others' expected contribution even in one-shot settings. Conditional cooperators are heterogeneous in the extent to which they match others' contributions. Many are 'imperfect' reciprocators in that they contribute slightly below what they expect others to contribute on average. In an experiment specifically designed to identify preference types, Fischbacher et al. (2001) classify 50% of their subjects as (partial) conditional cooperators and 30% as free-riders.² Others have roughly replicated this distribution of types in different subject pools (Kocher et al. 2008; Herrmann and Thöni 2009; Kamei 2012; Thöni et al. 2012). Conditionally cooperative behavior is consistent with a prosocial orientation (Van Lange 1999).

In repeated PGGs, conditional cooperators adapt their expectation of others' contribution on the basis of their experience of the average group contribution in the previous rounds (Fischbacher and Gächter 2010). The more free-riders and imperfect conditional cooperators there are, the lower group contribution will be. Conditional cooperators decrease their contribution accordingly, which causes the average to further decline. This explains the decrease of cooperation over time.

2.3 The PGG with sanctions

Sanctioning can be modeled by adding a second stage to the standard PGG. After all actors i have determined their contribution and observed the contributions of the other group members, they decide for every other group member j whether to pay an amount to punish and/or reward this actor. Let s_{ij} denote the amount actor i uses to sanction actor j . We assume here that an actor can only choose

² Virtually all remaining subjects were characterized as 'triangle' contributors (Fischbacher et al. 2001). These actors fully reciprocate others' expected contribution at 50% of the endowment, but their contribution declines when they expect others to contribute either more or less than this threshold.

whether or not to sanction, but not the magnitude of the sanction: s_{ij} is either a fixed amount $f > 0$ or zero. When the amount is used for punishment, a multiple k of f is subtracted from the payoff actor j obtained in the PGG. The same amount is added to the payoff of actor j when s_{ij} is used for reward. Thus, in addition to the payoff from the standard PGG, every actor j loses a total amount $k \sum_i s_{ij}$ of received punishment from all other actors i or gains this amount of received rewards. Moreover, every actor i forfeits $\sum_j s_{ij}$ by assigning sanctions to other actors j . This captures the essential features of how sanctions are executed in the PGG, denoted here as an IDR.³

In sanctioning systems with a CDR, all actors i likewise decide whether to pay an amount to sanction others. Sanctioning under a CDR is different from an IDR in the sense that a sanction is only implemented when at least a proportion p of all group members, save the prospective recipient, sanctions the same actor. Because we fixed the sanctioning amount s_{ij} to 0 or f , this implies that sanctions under a CDR are more severe than those under an IDR assigned by a smaller number of actors. Thus, every actor j loses an amount $k \sum_i s_{ij}$ in a punishment system with a CDR if the proportion q_j of actors i for whom $s_{ij} = f$ is larger than or equal to p . The same applies to the amount gained under rewards. If $q_j < p$ no sanction is executed, that is, actor j does not gain or lose money due to received sanctions. Moreover, the actors who proposed to sanction actor j do not pay the cost of sanctioning if $q_j < p$. Thus, every actor i who sanctions j loses an amount $\sum_{j: q_j \geq p} s_{ij}$.

We assume one-shot interactions.⁴ Thus, actors cannot benefit from group members who increase their contribution in subsequent games after being sanctioned. This implies that long-term incentives for sanctioning, which differ between IDRs and CDRs, are ruled out. Under these assumptions, rational selfish actors do not use costly sanctions regardless of the sanctioning of others. The Nash equilibrium of the one-shot PGG with sanctions under standard assumptions of rationality and selfishness is no sanctioning and no contributions. Although repeated interactions with sanction opportunities might be more realistic for many applications, we do stick to one-shot interactions also because in repeated interactions actors have alternative sanctioning mechanisms, too. For example, actors can reciprocate others' low contribution decisions by own low contribution

³ Note that details of this procedure can vary. For example, in many studies the amount s_{ij} used to sanction can be chosen freely by actors between 0 and some positive value.

⁴ In our experiment we employ random matching. Although subjects are likely to interact multiple times, they are not informed on the identity of others. It is common in the experimental literature to treat this as series of one-shot games. However, Botelho et al. (2009) show that, under random matching, subjects behave slightly different from subjects who play perfect stranger matching. The main difference is that subjects contribute zero more often under random matching, although they do not contribute less on average. We do not expect that these slight differences affect the difference between the experimental conditions we consider.

decisions in future interactions. This would lead to possible confounding effects of the exogenous sanctioning mechanisms we want to study with the endogenous sanctioning opportunities due to repeated interactions (cf. Fehr and Gächter 2002).

Given that not all actors are rational and selfish, some might sanction despite the prediction that follows from the assumptions of selfish rationality. Non-executed sanctions are costless, and group members are not informed about sanctions that were proposed by others but were not executed. Thus, non-executed sanctions cannot influence behavior of other actors than the ones who proposed the sanction. Therefore, actors have no incentive to take the probability that the sanction is executed into account when deciding whether or not to sanction under a CDR.⁵ Given these characteristics of the interaction situation, there is no reason to assume that actors make different *sanctioning* decisions under IDRs and CDRs.

We proceed with a review of empirical evidence and a theoretical account of contributing and sanctioning behavior in the PGG with an IDR. This reveals which actors allocate sanctions, and which behaviors are more likely to be sanctioned. Multiple individual sanctions for a given behavior imply a high consensus. Thus, given that the decision rule will not directly influence sanctioning decisions, those behaviors that are sanctioned individually by many are more likely to be sanctioned when a CDR is used. Behavior in the PGG with an IDR then allows to predict the likelihood that sanctions will be implemented under a CDR.

2.4 Behavior in the PGG with sanctions under an IDR

Despite the equilibrium prediction, empirical evidence shows that actors frequently use punishment under an IDR in one-shot settings. It is consistently found that punishment is assigned in accordance with enforcing cooperation. That is, actors receive more punishment the less they contribute (e.g. Carpenter and Matthews 2009; Casari and Luini 2009), and the less they contribute compared to the average contribution of the group (e.g. Fehr and Gächter 2000, 2002; Ones and Putterman 2007; Sefton et al. 2007; Carpenter and Matthews 2009; Ertan et al. 2009). This punishment is mostly executed by high contributors (e.g. Fehr and Gächter 2002; Sefton et al. 2007). It is also observed, however, that low contributors occasionally punish above-average contributors. This ‘perverse’ punishment is usually carried out by a small number of actors (Casari and Luini 2009). The extent to which it occurs varies greatly between subject pools, up to 50% of total punishment expenditure (Herrmann et al. 2008), but is typically estimated between 5% and 25% (Ostrom et al. 1992; Cinyabuguma et al. 2006; Ones and Putterman 2007;

⁵ When a sanction under a CDR is implemented, there are by definition multiple sanctioners. Thus, under a CDR actors know that implemented sanctions are severe for the recipient. Therefore, they might for example be reluctant to sanction mildly deviant behavior, for which a severe sanction might not be deserved. On the other hand, the chance that the sanction is implemented will also be smaller for light deviations. Thus, actors expect a more severe sanction with a lower probability of implementation the more are required to agree. There is no reason to assume this influences their sanctioning decisions.

Casari and Luini 2009; Ertan et al. 2009). The effect of cooperation-enforcing and perverse punishment differs. Below-average contributors increase their contribution in the subsequent round after being punished (e.g. Fehr and Gächter 2002), but for above-average contributors empirical evidence is mixed. Some studies show that above-average contributors decrease their contribution after being sanctioned (Masclet et al. 2003; Bochet et al. 2006; Ones and Putterman 2007); others find no effect of perverse punishment on contribution (Denant-Boemont et al. 2007; see also Ellingsen et al. 2012).

Like punishments, rewards are typically used to enforce cooperation in one-shot settings. High contributors tend to reward other high contributors (Walker and Halloran 2004; Sefton et al. 2007; Sutter et al. 2010; Ellingsen et al. 2012; Choi and Ahn 2013). However, while rewards are mainly allocated to above-average contributors, it is often less clear than for punishment that the amount of rewards received increases with the (positive) deviation from the average group contribution (Walker and Halloran 2004; Sefton et al. 2007; Nosenzo and Sefton 2012; Choi and Ahn 2013; but see Ellingsen et al. 2012). Also, in repeated PGGs where actors can identify each other it is found that rewards are frequently used in every successive interaction (Rand et al. 2009; Milinski and Rockenbach 2011; Ellingsen et al. 2012), while the use of rewards declines over time in fixed groups when actors cannot infer who rewarded them (Sefton et al. 2007; Choi and Ahn 2013). As with punishments, the effect of rewards differs with the recipients' contribution. Above-average contributors are found to contribute more in the subsequent interaction the more rewards they receive, while below-average contributors decrease their contribution the more they are rewarded (Ellingsen et al. 2012).

In repeated interactions in fixed groups, contributions under reward are sometimes found to be lower than those under punishment (Sutter et al. 2010 low leverage; Milinski and Rockenbach 2011; Wiedemann et al. 2011; Drouvelis and Jamison 2012; Nosenzo and Sefton 2012) although others did not find a difference, at least until the final periods (Sefton et al. 2007; Rand et al. 2009; Sutter et al. 2010 high leverage; also see Balliet et al. 2011; Choi and Ahn 2013). However, in repeated one-shot settings, which are most similar to our experiment, it is found that contributions are lower under rewards than under punishment (Choi and Ahn 2013).

2.5 Non-selfish utility in the PGG with sanctions

Rational selfish free-riders never sanction when this is costly. However, anticipation on being sanctioned will induce them to contribute, provided that the loss due to received punishment or gain from rewards offsets the payoff advantage of free-riding (Fehr and Fischbacher 2004). Non-selfish actors could derive utility from sanctioning defectors even in one-shot interactions (Diekmann and Voss 2003). These cooperation-enforcing punishers are sometimes classified as a separate type of actor, which partly but not completely overlaps with

conditional cooperators in the PGG without punishment (e.g. Ostrom 2000; Ones and Putterman 2007).

Empirical evidence is indeed consistent with the assumption that people derive utility from punishing and rewarding in one-shot settings. Fehr and Gächter (2002) already noted that subjects experience anger when they observe free-riding in a hypothetical situation. This anger increases the more the free-rider deviates from the average contribution of others. Casari and Luini (2009) show that punishment decisions are not influenced by information that others already punished the recipient. Thus, subjects do not care so much about actors being punished, but derive utility from the act of punishing. Fudenberg and Phatak (2010) show that subjects punish even when the recipient is not informed on the punishment, implying that punishment cannot influence future cooperation. In a neurobiological experiment, De Quervain et al. (2004) show that the human reward system is activated in the brain of an actor punishing a defector. Utility from rewarding is addressed by Dawes et al. (2007), who conduct an experiment in which subjects can decide on a costly in- or decrease of a random amount of tokens other subjects had received. They find that subjects who afterwards indicate more anger and annoyance towards those with a high amount also spend more to increase low and reduce high amounts received by others. Yet, despite utility derived from sanctioning, it is found that actors sanction less the higher the costs of sanctioning are (Anderson and Putterman 2006; Carpenter 2007; Nikiforakis and Normann 2008; Vyrastekova and Van Soest 2008; Sutter et al. 2010). Thus, actors take their own payoff into account in sanctioning decisions (Fehr and Fischbacher 2004).

As mentioned above, some actors use sanctions perversely. Although they are relatively rare, perverse sanctioners constitute a separate type of actors. These actors free-ride in the PGG, and subsequently punish high contributors (e.g. Cinyabuguma et al. 2006; Herrmann et al. 2008; Gächter and Herrmann 2009; Chaudhuri 2011). A motive for perverse punishment might be revenge on previous punishment received from high contributors (Ostrom et al. 1992; Fehr and Gächter 2000; Denant-Boemont et al. 2007; Nikiforakis 2008), a desire to increase relative payoff advantage of free-riding (Fehr and Gächter 2000), or a dislike of do-gooders or norm violators (Monin 2007; Ones and Putterman 2007; Gächter and Herrmann 2009). Alternatively, it could be that actors occasionally punish high contributors by mistake (Fehr and Gächter 2000). Rand et al. (2010) and Rand and Nowak (2011) show that punishment of cooperators can be evolutionary stable, thus providing a potential explanation for the fact that perverse punishment can drive out cooperation. Perverse rewards, i.e. rewards targeted at free-riders, just as perverse punishments, increase the payoff discrepancy between high and low contributors. Hence, they are potentially equally detrimental for cooperation (Ellingsen et al. 2012).

Punishment and reward are used in different ways. The possibility of being punished might be enough to deter free-riding, such that there is no need to

actually allocate punishment. However, when an actor makes a high contribution, rewards actually have to be carried out sufficiently often to induce free-riders to contribute (Dari-Mattiacci and De Geest 2010). Thus, when contributions in a population increase due to the existence of a sanctioning system, more rewards than punishments have to be allocated. In one-shot settings, actors cannot establish a norm of direct mutual rewarding. They are therefore unsure whether the costs of allocating rewards will be offset by reciprocation (Rand et al. 2009). This makes rewards more expensive than punishments in the one-shot PGG. As stated above, more expensive sanctioning implies that less sanctions are assigned. This explains why, without opportunities for directly reciprocating received rewards, actors initially attempt to reward but eventually give up when others do not continue to reward as well.

2.6 Micro-level hypotheses

Before turning to differences in contribution levels between IDRs and CDRs, we capture the framework developed for micro-motives, that is, contributing and sanctioning behavior of individual actors, in a number of hypotheses. These hypotheses are based on empirical regularities observed in previous experiments. The hypotheses will be used as a micro-level framework summarizing which actors are likely to sanction, and how actors react to receiving cooperation-enforcing or perverse sanctions. When theorizing about the effect of sanctioning decision rules on contributions, we assume that actors behave as summarized in this framework.

We first derive hypotheses on sanctioning behavior. Although perverse punishment is sometimes observed, punishment is usually allocated by cooperation-enforcing high contributors. Accordingly, we hypothesize that actors are more likely to punish others the more they contributed themselves.

Hypothesis 1: The more an actor contributes, the higher this actor's likelihood to assign punishment.

Punishment of high contributors is more often targeted at free-riders than punishment of low contributors, who might punish perversely. Thus, the more an actor contributed the more likely he is to punish a free-rider. This implies that we expect an interaction between the contribution of the actor allocating punishment and the contribution of the recipient on the likelihood to sanction. We argue that actors perceive free-riding both in the sense of the recipient contributing a low amount and in the sense of contributing less than the other group members. This means that low as well as below-average contributors are likely to be punished by high contributors.

Hypothesis 2a: The more an actor contributes, the more this actor's likelihood of assigning punishment decreases with the contribution of the recipient.

Hypothesis 2b: The more an actor contributes, the more this actor's likelihood of assigning punishment increases with the negative deviation of the recipient from the group average contribution.

Also reward is predominantly allocated by high contributors.

Hypothesis 3: The more an actor contributes, the higher this actor's likelihood to assign reward.

High contributors are more likely to reward other high contributors. This applies both in an absolute sense, and compared to the average of other group members. Again, we hypothesize an interaction between the contribution of the rewarding actor and the contribution of the recipient.

Hypothesis 4a: The more an actor contributes, the more this actor's likelihood of assigning reward increases with the contribution of the recipient.

Hypothesis 4b: The more an actor contributes, the more this actor's likelihood of assigning reward increases with the positive deviation of the recipient from the group average contribution.

Unlike punishments, in order to enforce cooperation rewards have to be allocated repeatedly to high contributors. They are therefore costly to maintain when direct reciprocation is impossible. Accordingly, the likelihood of rewarding decreases over rounds.

Hypothesis 5: The more rounds have already been played, the lower the likelihood that rewards are allocated.

We now turn to the effect of sanctions on contribution. Receiving punishment leads to conformation to behavior of other actors, in order to avoid receiving punishment in future interactions. Free-riders thus increase and high contributors decrease contribution the more they are punished. Consequently, their contribution is more in line with others' average.

Hypothesis 6: The more an actor contributing below the average is punished, the more this actor contributes in the subsequent interaction.

Hypothesis 7: The more an actor contributing above the average is punished, the less this actor contributes in the subsequent interaction.

Rewards strengthen current deviations from average behavior. Above-average contributors will thus contribute more and below-average contributors less the more they are rewarded, provided they did not already contribute the full endowment or free-ride completely, respectively.

Hypothesis 8: The more an actor contributing above the average is rewarded, the more this actor contributes in the subsequent interaction.

Hypothesis 9: The more an actor contributing below the average is rewarded, the less this actor contributes in the subsequent interaction.

2.7 Macro-level effects of CDRs

Only the sanctions on which required consensus is reached are executed under a CDR. Given sanctioning behavior as predicted in the micro-level hypotheses, it is likely that there will be more consensus on some sanctions than on others. This gives rise to different contribution levels under IDRs versus CDRs. Macro-level hypotheses differ for punishment and reward.

Under an IDR, all allocated punishments are carried out. This implies that high contributors will frequently punish free-riders. Free-riders will receive more punishment the less they contribute in absolute sense and compared to the others. Also, perverse punishers have the opportunity to punish high contributors.

The situation is different when only those sanctions are implemented to which a majority of actors consents. A large proportion of actors derives utility from sanctioning. It is therefore likely that majority consent is often reached on punishment of free-riders. The more a free-rider deviates from the average, the higher the chance that consent is reached. Conversely, when perverse punishment is relatively rare, as is typically found, it will be unlikely that a majority of actors agrees on punishing a high contributor. Thus, a majority sanctioning system will mitigate perverse punishment while at the same time cooperation-enforcing punishment is likely to be implemented. We therefore expect a majority decision rule to lead to higher contribution levels than an IDR.

Hypothesis 10a: Contribution is higher under a majority than under an individual punishment decision rule.

Some previous studies indeed found that majority consent is sufficient to rule out perverse punishment, but that cooperation-enforcing punishment could still be implemented. Casari and Luini (2009) found that punishment was more effective when two out of four actors had to agree on sanctioning a fifth. Perverse punishment was to a large extent ruled out under this decision rule. Likewise, Ertan et al. (2009) let subjects choose whether or not to enable punishment of high contributors. While this was sometimes favored by a number of free-riders, it was never implemented because a majority opposed the possibility.

Under a unanimity decision rule punishment is only executed when all remaining group members decide to punish an actor. Perverse punishment is therefore even less likely than under a majority decision rule. However, also for cooperation-enforcing punishment a unanimity decision rule requires a very high proportion of actors willing to punish. Therefore, it will be difficult to implement any punishment at all. Conversely, under an IDR there could be perverse punishment, although the vast majority of punishment should be targeted at below-average contributors. It is therefore likely that contribution levels under a unanimity punishment decision rule are lower than under an individual rule.

Hypothesis 10b: Contribution is higher under an individual than under a unanimity punishment decision rule.

As explained above, continuous need of rewarding makes reciprocating through rewards more expensive than through punishment, which causes the use of rewards to decline (Dari-Mattiacci and De Geest 2010). Thus, more punishment than reward will be executed under every decision rule, making sanctioning through punishment more effective. Therefore, we expect that for every decision rule contribution is higher under punishment than under reward.

Hypothesis 11: For every decision rule, contribution is higher under punishment than under reward.

The more actors are required for a reward to be executed, the more likely it is that too many actors give up on using rewards. Thus, the more actors are required the more likely it is that consensus cannot be reached anymore. Also, perverse rewards have to be carried out when an actor free-rides in anticipation on being rewarded. Perverse rewards are thus likewise costly to maintain. Therefore, while perverse rewards might be occasionally allocated it is unlikely that they are persistently problematic for enforcing cooperation. Thus, rewards under an IDR are not thwarted by perverse sanctions as much as punishment, while it is difficult to raise enough actors to agree on rewards under a CDR. The more actors are required to agree, the more problematic enforcing cooperation becomes. Accordingly, we hypothesize that the more actors are required to agree on rewards, the less rewards will be carried out and the lower contribution levels are. Thus, the macro-level hypotheses on rewards are partly different from those on punishment.

Hypothesis 12a: Contribution is higher under an individual than under a majority rewarding decision rule.

Hypothesis 12b: Contribution is higher under a majority than under a unanimity rewarding decision rule.

3. Experimental design

In the experiment, subjects participated in interaction situations based on the PGG as described above with group size $n=4$; endowment $w=20$, and multiplier $m=1.6$. The outcome of the game represented points that subjects earned. After the experiment, subjects received 1 eurocent for every 60 points earned.

The experiment comprised three parts. In the first part, preferences for conditional cooperation were assessed using a measure designed by Fischbacher et al. (2001). First, subjects decided on an unconditional contribution, i.e. how much to contribute in the PGG in a group with three other subjects. Second, subjects made this same decision conditional on others' average contribution. Thus, they decided how much they would contribute for every possible average of the three other group members (strategy method, Selten 1967). The more conditionally cooperative a subject is, the more contribution should increase with others' average.

Subjects were randomly matched in groups of four. For three randomly chosen group members, payoff was calculated based on the unconditional contribution. For the fourth group member the conditional contribution corresponding to the average unconditional contribution of the three others was used. This makes both decisions incentive-compatible. Note that conditionally cooperative preferences were always assessed at the beginning of a session, prior to playing the actual PGGs. Fischbacher and Gächter (2010) measured conditional cooperation using a similar design, administered either at the start or end of the experiment. They did not find a sequence effect, suggesting that measuring preferences does not significantly influence subsequent behavior.

In the second part of the experiment, the standard PGG as described above was played for 10 rounds. Between the rounds, subjects were randomly rematched into different groups. They could not infer their group members' previous decisions. After every round, subjects were informed about the contribution of the others in their group and their own payoff. Numerous previous experiments have administered baseline games before the experimental treatments (cf. Sefton et al. 2007; Casari and Luini 2009). No treatment effects were found in experiments where the order of baseline and punishment treatments was randomized (e.g. Fehr and Gächter 2002; Herrmann et al. 2008).

In the third part, the PGG with sanctions was employed. In every session, 10 rounds were played with only punishment and 10 rounds with only reward; the order varied between sessions. Both reward and punishment took place in one of three experimental conditions; individual, majority, or unanimity. In all three conditions, subjects first decided upon a contribution. Subsequently, they were informed about contributions of their group members and decided for all three others separately whether to sanction this person. If executed, a sanction added or subtracted six points from the earnings of the recipient at a cost of two points. This cost ratio of 1:3 is often used in PGG experiments (cf. Fehr and Gächter 2002). The effect and cost of the sanction were chosen to ensure that receiving a sanction has a severe impact on payoffs. Because the amount by which actors could sanction was fixed, the severity of the sanction is equal to the number of actors sanctioning.⁶

In the individual condition, all assigned rewards and punishments were implemented. Subjects who received multiple sanctions were sanctioned by the cumulative amount while all subjects allocating the sanction paid the cost of two points. The procedure in the majority condition was exactly the same, except that the sanction was only executed when at least two group members wanted to sanction the same recipient. Thus, an actor sanctioned by two others lost 12 points, while both sanctioning actors lost 2 points. In the unanimity condition, the sanction was only executed when it was requested by all three remaining group members. When the number of subjects who wanted to sanction was insufficient in the majority or

⁶ The 1:3 reward ratio enables increasing group earnings through mutual rewarding. However, note that our random matching scheme excludes direct reciprocity. Subjects are therefore unlikely to unilaterally reward all others for the purpose of increasing efficiency.

unanimity condition, the sanction was not executed and no costs had to be paid. Note that the labels “majority” and “unanimity” imply that a subject is not involved in the decision of sanctioning him- or herself. Thus, only the three other subjects determine whether the fourth subject is going to be sanctioned. After each round, subjects were informed about all sanctions that had been executed in their group but could not infer who allocated them. No information was provided about sanctions that were not executed. Again, subjects were randomly rematched between the rounds.

The experiment was programmed using z-Tree (Fischbacher 2007) and conducted at the ELSE laboratory of Utrecht University. Subjects were recruited using the online recruiting system ORSEE (Greiner 2004). Twelve sessions were held, four in each experimental condition of which two with reward first and two with punishment first. Instructions were provided on paper. It was made clear that the instructions were always truthful and identical for all subjects in a session. In the first set of instructions, the standard PGG and the first two parts of the experiment were explained. It was announced that there would be further tasks, but not what these tasks entailed. These instructions included a number of control questions, which appeared on the computer screen. When a subject did not answer correctly to a question, the answer was explained on the screen. Additional instructions, adapted for each experimental condition, were provided for the reward as well as for the punishment part. The options in the PGG were labeled in a neutral way: punishment and reward were called ‘subtracting’ and ‘adding’ points, respectively.

A total number of 184 student subjects participated in the experiment (32% male; 34% economics major). Both the majority and unanimity sessions comprised 64 subjects in total, while 56 subjects were in a session which was held in the individual condition. Payoffs averaged €12.50, with a minimum of €8.50 and a maximum of €15.

4. Method and results

4.1 Descriptive results

All subjects participated first in the baseline, and subsequently in reward as well as punishment of one of the conditions. A Mann-Whitney test revealed no significant effect of the order in which punishment and reward treatments were administered on average contribution in either the reward ($z=1.601$; $p=0.11$) or punishment ($z=1.441$; $p=0.15$) games.⁷ However, since these p -values are relatively low we check the robustness of our parametric analyses, in which we combine the two sanctioning treatments, against analyses in which only the first sanctioning treatment is included.

Figure 1 shows the average contributions in the PGGs over the rounds in the baseline and in each experimental condition. Contributions are initially around

⁷ Reported p -values of all non-parametric tests are two-sided. The experimental sessions are used as independent observations.

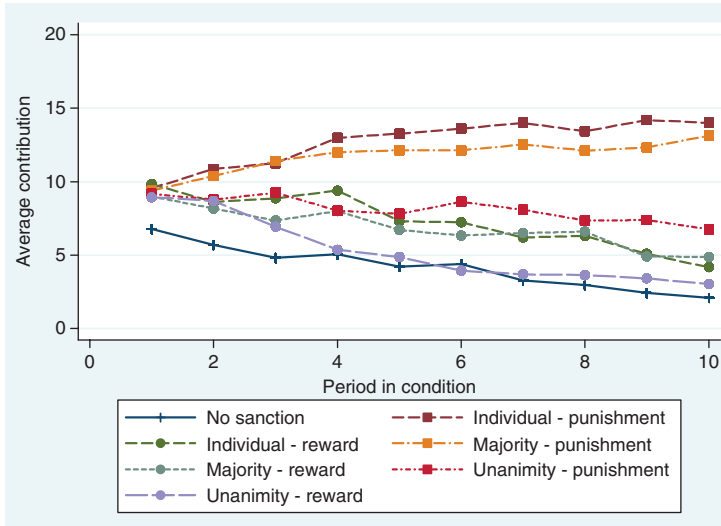


Figure 1: Average contribution in the PGGs, separated for each round and experimental condition.

50% of the endowment. This is in line with previous findings (Ledyard 1995). After the first round, Figure 1 shows strong differences in contribution levels between the conditions. Contributions in the baseline decline to almost zero. Conversely, individual and majority punishment are the only conditions under which contributions increase over time. A Wilcoxon signed rank test confirms that average contribution is higher in the reward than the baseline ($z=2.432$; $p=0.02$) and in the punishment than the reward conditions ($z=3.059$; $p<0.01$). For both reward and punishment the individual and majority conditions lead to higher contributions than unanimity, although only the difference between individual and unanimity punishment is significant in a Mann-Whitney test ($z=2.309$; $p=0.02$).

Overall average profits are higher in the reward than in both the punishment and baseline treatments (Wilcoxon signed rank test: $z=2.589$; $p=0.01$ for baseline vs. reward; $z=2.981$; $p<0.01$ for punishment vs. reward; $z=1.098$; $p=0.27$ for baseline vs. punishment). However, this is related to our reward technology, which enables earnings to be higher in the reward than the other treatments. Highest possible group earnings are achieved with full contribution in baseline and punishment, and with full contribution and mutual rewarding in the reward treatments. When we consider average earnings as a proportion of the highest possible, this proportion is higher in both punishment ($z=3.059$; $p<0.01$) and baseline ($z=3.059$; $p<0.01$) than in the reward treatments.

When a subject was punished in the majority condition, in 58% of the cases this was by one person only and therefore the punishment was not

carried out. Likewise, in 81% of the cases in which a subject was punished in the unanimity condition the required number of three sanctioning subjects was not reached. For reward, in 72% of the cases in which someone was rewarded in the majority condition and in 97% of the cases under unanimity the reward was not implemented. In line with previous research, 25% of punishments were targeted at subjects contributing the average of other group members or more. Of these, 91% and 98% were not implemented in majority and unanimity, respectively. 33% of rewards were targeted at below-average contributors, of which 89% and 100% were not implemented under majority and unanimity.

Figure 2 shows the average number of sanctions allocated and average number of sanctions carried out for different deviations of the recipient from the average contribution of the other group members. Note that between one and three other group members can propose to sanction. Figure 2 shows a clear trend of more punishment proposed on average the more the recipient negatively deviates from the average contribution of others. Also, more rewards are proposed for above-average contributors, but it is not so clear that more rewards are proposed the further the deviation.

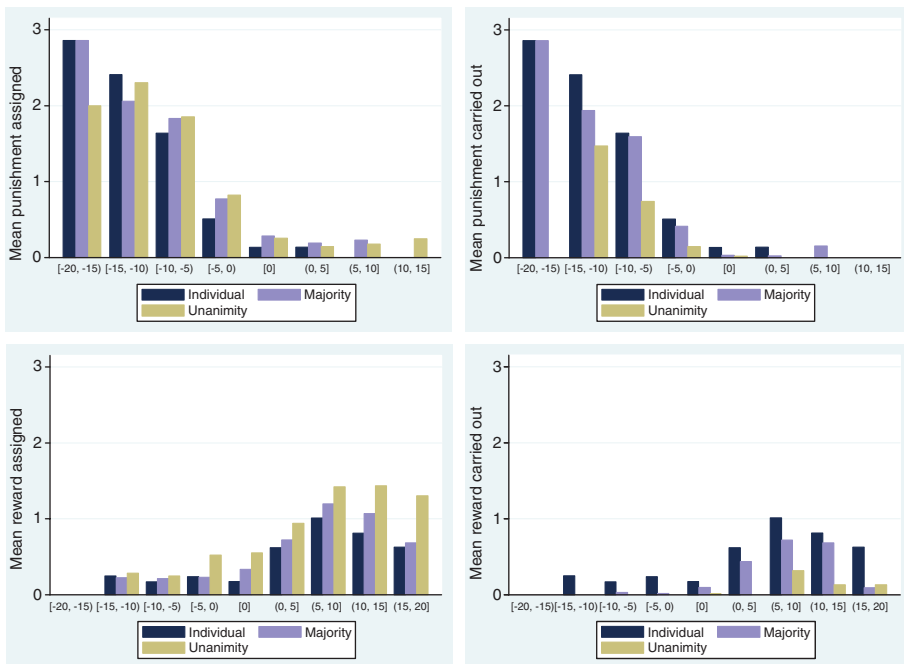


Figure 2: Average punishment (above) and reward (below) assigned (left) and carried out (right) for different deviations from the average contribution of other group members, separated for each experimental condition.

4.2 Contribution – methods

The first dependent variable, contribution, is measured as the contribution decisions of subjects in the PGG. First, we test macro-level hypotheses by comparing dummies for the experimental conditions individual, majority, and unanimity punishment and reward. These are less conservative tests for the differences between conditions than the comparisons in the previous subsection, because the interdependencies between the observations are modeled in more detail. Still, the results mainly reconfirm the differences that resulted from the non-parametric tests. Second, we test the micro-level hypotheses *explaining* differences between experimental conditions based on individual decision patterns. Punishment and reward conditions are analyzed separately.

In the micro-level models, sanctions received are measured as the number of others who had sanctioned the subject in the previous round. Only executed sanctions are included. Furthermore, three dichotomous variables indicate whether in the previous round a subject had contributed more than 4 points below the average of other group members, more than 4 points above the average, or did not deviate from the average by more than 4 points. These three dummies for previous deviation are interacted with the number of sanctions received to test whether the effect of being sanctioned is different for above- and below-average contributors.

Previous deviation was measured using dummies for more than 4 points higher/lower rather than a continuous variable indicating the precise extent of the deviation. This is because a continuous variable interacted with received reward tests if subjects increase (decrease) their contribution more, the higher (lower) the contribution for which they were rewarded. This is unrealistic, since contribution is limited between 0 and 20. The boundaries of 4 points from the average are chosen such that the deviation is substantial enough for subjects to perceive sanctions as clearly norm-enforcing or perverse. Accordingly, log-likelihoods of models with different boundaries are equal to or lower than those of the models presented here. We control for the subjects' contribution in the previous round, round number, treatment order, and experimental condition. Furthermore, preference for conditional cooperation is included, measured as the slope of the conditional contribution assessed in the first part of the experiment. The steeper the slope, the more a subject indicated to contribute more when others do so as well.⁸

We use Tobit regression to take into account that contribution has a limited range, between 0 and 20, of which both extremes are often chosen. The units of analysis are decisions in the PGGs. Random effects at the subject level are included to model that decisions are nested in subjects, since every subject makes multiple contribution decisions. Also, within a session subjects often encounter

⁸ Two subjects whose slopes are zero, but who do make positive conditional contributions (both unconditional contributors) are excluded from the analysis. A zero slope thus indicates a preference for unconditional free-riding. Note that excluding these subjects did not influence the results.

others with whom they or their group members have interacted previously. Thus, subjects are interdependent within sessions. It is not possible to include both the subject and session level in a three-level Tobit model. Therefore, all models were replicated using multilevel linear regression, in which both subject and session level random effects are included but where contribution is treated as if its range is unlimited. Also, we estimated the models using Tobit regression with random effects at the session level to test if disregarding this level in the models presented below influenced the results, and we ran a Tobit model with robust standard errors adjusted for clustering within sessions. The latter model provides the most conservative way of correcting for the clustering of observations and, therefore, might underestimate the significance of some effects. Given the limited effect of the session level in, e.g. the three-level linear regression model, we have considerable confidence in the estimations of the two-level Tobit models with random effects for subjects reported in the tables. Finally, we examined the possible effects of punishment and reward treatment order in more detail by rerunning all models with only the first treatment that subjects participated in included. Effects of treatment order and robustness of the results in alternative analyses are discussed for every model separately below.

4.3 Contribution – results

Table 1 shows differences in contribution decisions between the experimental conditions. The baseline condition, in which every subject participated, serves as a reference. Contributions in all experimental conditions except unanimity reward were higher than in the baseline, although the effect of majority reward is insignificant when we adjust for clustering within sessions. Contrary to Hypothesis 10a, contribution under punishment is higher in the individual than

Table 1: Tobit regression on contribution decisions with random effects at subject level (5460 decisions, of which 2376 censored, by 182 subjects).

	Model 1	
	Coeff.	S.D.
Baseline	Ref.	
Punishment – Individual	13.938**	0.518
Punishment – Majority	10.239**	0.464
Punishment – Unanimity	5.866**	0.479
Reward – Individual	6.184**	0.528
Reward – Majority	2.770**	0.474
Reward – Unanimity	0.340	0.501
Constant	0.786	0.522
σ_u	6.372**	0.380
σ_e	7.934**	0.113
Log Likelihood	-12773.784	

*Significant at .05-level; ** Significant at .01-level (2-sided).

the majority condition ($\chi^2(1)=29.51$; $p<0.01$). The other macro-level hypotheses are confirmed. Contribution under punishment is higher in the individual than the unanimity condition ($\chi^2(1)=136.58$; $p<0.01$), confirming Hypothesis 10b. As predicted in Hypothesis 11, contribution is higher under punishment than reward in the individual ($\chi^2(1)=228.83$; $p<0.01$), majority ($\chi^2(1)=246.01$; $p<0.01$), and unanimity ($\chi^2(1)=122.01$; $p<0.01$) condition. Finally, contribution under reward is higher in the individual than the majority condition ($\chi^2(1)=23.76$; $p<0.01$) and higher in the majority than the unanimity condition ($\chi^2(1)=12.79$; $p<0.01$). This confirms Hypotheses 12a and 12b. All differences between decision rules are insignificant in the conservative model that accounts for clustering in sessions, but remain highly significant in other model specifications. The differences between punishment and reward remain significant in every alternative specification.

Because we want to exclude that the support for the hypotheses confounds with effects of subjects playing a punishment and reward treatment after each other, we also consider effects of the ordering of treatments. Contributions in the punishment conditions are lower when punishment was the first compared to when it was the second treatment. In the reward conditions, contributions are higher when it was the first treatment. Still, when we only consider the first treatments subjects participated in, contributions are higher in individual than in majority conditions, although this difference becomes insignificant for reward ($\chi^2(1)=0.93$; $p=0.34$). Also, contributions are higher in majority than in unanimity conditions. Finally, contributions are higher in the individual and unanimity punishment conditions than in the related reward conditions. Only in the majority condition this difference disappears ($\chi^2(1)=0.16$; $p=0.69$). Hence, the confirmation of this part of Hypothesis 11 should be interpreted with caution.

The micro-level model for the punishment conditions is presented in Table 2. Only main effects are included in Model 2. Several control variables are significant. Contribution is lower in the unanimity compared to the individual condition and when punishment was administered first, and higher the more a subject contributed in the previous round. The difference between the individual and majority condition is not significant in this model. Subjects who contributed 4 points or more below the average increase and subjects who contributed above the average decrease their contribution compared to around-average contributors. Also, contribution is higher the more punishment was received previously.

Interaction effects are included in Model 3. The main effect of punishment is excluded from this model, so the three interactions represent the effect of received punishment for the three groups of subjects belonging to specific deviations from the mean contribution. The model shows that subjects contributing below the average increase their contribution more, the more they are punished. Hypothesis 6 is thus confirmed. The insignificant main effect of negative deviation indicates that subjects who contributed below the average but were not punished do not significantly increase their contribution compared to around-average contributors. Subjects who contributed above the average decreased their contribution if they

Table 2: Tobit regression on contribution decisions in the punishment conditions with random effects at subject level (1638 decisions, of which 345 censored, by 182 subjects).

	Exp. direction	Hyp. nr.	Model 2		Model 3	
			Coeff.	S.D.	Coeff.	S.D.
Previous punishment received			1.006**	0.170		
Prev. neg. deviation > 4			1.102*	0.440	0.045	0.514
× punishment received	+	6			1.598**	0.226
Prev. deviation ≤ 4			Ref.		Ref.	
× punishment received					0.285	0.256
Prev. pos. deviation > 4			-1.905**	0.347	-1.985**	0.353
× punishment received	-	7			0.361	0.829
Previous contribution			0.642**	0.044	0.653**	0.044
Slope conditional contribution			0.986	0.559	0.898	0.543
Period			-0.014	0.043	-0.015	0.043
Individual			Ref.		Ref.	
Majority			-0.057	0.614	-0.101	0.596
Unanimity			-2.907**	0.650	-2.838**	0.633
Punishment treatment first			-1.328**	0.506	-1.248*	0.492
Constant			10.568**	0.594	10.199**	0.591
σ_u			2.964**	0.263	2.853**	0.262
σ_e			4.189**	0.093	4.180**	0.093
Log Likelihood			-4130.177		-4122.261**	

*Significant at .05-level; ** Significant at .01-level (2-sided).

had not been punished, but did not decrease their contribution further after receiving punishment.

Thus, no support is found for Hypothesis 7. This might be due to the relatively limited amount of sanctioning against high contributors even in the individual condition. The effect remains insignificant in a separate analysis of the individual condition.

These findings in Models 2 and 3 are similar in a multilevel model, with random effects and clustering at session level, and in a model in which only the first treatments are considered. All hypothesis-related effects are robust.

Model 4 in Table 3 shows the determinants of contribution decisions in the reward conditions. In this model the differences between experimental conditions and treatment order are not significant. The other control variables are significant; contribution is higher the more conditionally cooperative a subject is and the more a subject contributed previously, and decreases over rounds. Subjects who previously contributed above the average decrease and those who contributed below the average increase their contribution compared to around-average contributors. Finally, the more rewards a subject had previously received, the higher the contribution.

In Model 5, the interaction effects are included. Again, the three interactions represent the separate main effects. This shows that subjects who had contributed above the average significantly decrease their contribution. However, the decrease

Table 3: Tobit regression on contribution decisions in the reward conditions with random effects at subject level (1638 decisions, of which 981 censored, by 182 subjects).

	Exp. direction	Hyp. nr.	Model 4		Model 5	
			Coeff.	S.D.	Coeff.	S.D.
Previous reward received			1.571**	0.511		
Prev. neg. deviation > 4 × reward received	–	9	2.588**	0.869	2.491**	0.884
Prev. deviation ≤ 4 × reward received			Ref.		Ref.	
Prev. pos. deviation > 4 × reward received	+	8	–4.876**	0.994	–5.751**	1.083
Previous contribution			0.882**	0.091	0.901**	0.092
Slope conditional contribution			5.521**	1.563	5.507**	1.571
Period			–0.752**	0.123	–0.753**	0.123
Individual			Ref.		Ref.	
Majority			0.580	1.703	0.522	1.715
Unanimity			–2.669	1.746	–2.792	1.758
Reward first treatment			2.649	1.384	2.727	1.391
Constant			1.971	1.525	1.910	1.551
σ_u			8.135**	0.774	8.191**	0.775
σ_e			9.593**	0.316	9.560**	0.315
Log Likelihood			–3063.380		–3061.248	

*Significant at .05-level; ** Significant at .01-level (2-sided).

was significantly weaker the more they were rewarded. This confirms Hypothesis 8. Very few subjects received rewards after a below-average contribution, and virtually all were ruled out by majority and unanimity. Hence, we find no significant effect of being rewarded for around-average or below-average contributors. The effect remains insignificant in a separate analysis of the individual condition. Hypothesis 9 is not confirmed. Again, findings are similar in a multilevel model, with random effects and clustering at session level, and in a model in which only the first treatments are considered. All hypothesis-related effects are robust.

4.4 Sanctioning – methods

The second dependent variable in the analysis of the micro-level framework are the decisions whether or not to sanction. These are three observations for each subject in each period, one for every other group member.

The first independent variable is a subjects' own contribution. Furthermore, contribution of the recipient is included as a continuous variable. Deviation of the recipient from the average of others is measured as the contribution of the recipient minus the average of the other group members. The variable positive deviation includes all positive values of this measure, negative values are set to zero. Absolute negative deviation represents the extent of the deviation of all negative values, zero for positive deviations. For punishment, the contribution

and absolute *negative* deviation of the recipient are interacted with the subjects' own contribution to test whether high contributors are more likely to punish the less the recipient contributes, and the further he deviates from the average. For reward, contribution and *positive* deviation of the recipient are interacted with subjects' contribution. We control for experimental condition, treatment order, slope of the conditional contribution, and for sanctions assigned and received by the subject in the previous round.

We use logistic regression to analyze the dichotomous sanctioning decisions. Every subject makes three sanctioning decisions, one for every other group member, in all ten periods. Decisions are thus nested within periods and subjects. A multilevel intercept-only model with decisions nested in periods and subjects revealed that variance at the period level is negligible for both punishment and reward decisions. We therefore use multilevel models with decisions nested only in subjects. All models were repeated using only the first treatment subjects participated in. We discuss the treatment effects of all models below.

4.5 Sanctioning – results

Models on punishment decisions are displayed in Table 4. Model 6 shows that there are no differences between the experimental conditions in the likelihood that a subject decides to punish another. We do find that subjects who have received

Table 4: Multilevel logistic regression on decisions whether to punish nested in subjects (4914 decisions by 182 subjects).

	Exp. direction	Hyp. nr.	Model 6		Model 7	
			Coeff.	S.D.	Coeff.	S.D.
Contribution	+	1	0.091**	0.018	0.121**	0.019
Contribution recipient × Contribution	–	2a	–0.220**	0.027	–0.167**	0.029
Positive deviation recipient			–0.018	0.032	–0.075*	0.034
Absolute neg. deviation recipient × Contribution	+	2b	0.283**	0.029	0.211**	0.033
Round			0.029	0.020	0.050*	0.021
Individual			Ref.		Ref.	
Majority			0.203	0.394	0.239	0.444
Unanimity			–0.225	0.410	–0.188	0.458
Slope conditional contribution			0.295	0.357	0.145	0.402
Previous punishment received			0.301**	0.066	0.268**	0.067
Previous punishment assigned			0.291**	0.067	0.255**	0.068
Punishment first treatment			–0.604	0.323	–0.861*	0.364
Constant			–1.500**	0.371	–1.311**	0.409
σ_u			1.930**	0.168	2.187**	0.190
Log Likelihood			–1505.348		–1444.317**	

* Significant at .05-level; ** Significant at .01-level (2-sided).

† Hypothesized effect not significant when only the first treatment is considered.

or have allocated punishment in the previous round are more likely to punish. The likelihood of punishing increases with contribution, confirming Hypothesis 1. Also, the more a recipient negatively deviates from others' contribution, the higher the likelihood that punishment is allocated while no effect is found for positive deviation. Finally, the more a group member contributes, the less likely subjects are to punish this person.

Model 7 shows a significant interaction effect of contribution with the contribution of the recipient, confirming Hypothesis 2a. A significant interaction with negative deviation of the recipient confirms Hypothesis 2b. High contributors are thus more likely to punish the less a recipient contributes in absolute sense, and relative to the average of others.

The effect that high contributors punish especially others who contribute less than average (Hypothesis 2b) is not found if we only consider the first treatment for Model 7. This is probably due to the lower number of observations when only one treatment is included, which makes it more difficult to disentangle the different reasons why high contributors punish others.

Table 5 shows the models on reward. Main effects included in Model 8 show that subjects in the unanimity condition are more likely than in the individual condition to allocate rewards. Furthermore, subjects are more likely to reward the more rewards they had allocated in the previous period. The effect of period is

Table 5: Multilevel logistic regression on decisions whether to reward nested in subjects (4914 decisions by 182 subjects).

	Exp. direction	Hyp. nr.	Model 8		Model 9	
			Coeff.	S.D.	Coeff.	S.D.
Contribution	+	3	0.070**	0.012	0.037**	0.013 [†]
Contribution recipient			0.157**	0.019	0.150**	0.021
× Contribution	+	4a			0.005**	0.002 [†]
Positive deviation recipient			0.032	0.020	0.074**	0.023
× Contribution	+	4b			0.007**	0.003
Absolute neg. deviation recipient			-0.115**	0.025	-0.047	0.026
Round	-	5	-0.052*	0.024	-0.065**	0.024
Individual			Ref.		Ref.	
Majority			0.349	0.502	0.361	0.489
Unanimity			1.364**	0.506	1.311**	0.492
Slope conditional contribution			0.745	0.452	0.662	0.439
Previous reward received			-0.179	0.092	-0.187	0.100
Previous reward assigned			0.346**	0.075	0.350**	0.076
Reward treatment first			0.117	0.403	0.126	0.392
Constant			-3.567**	0.442	-3.517**	0.430
σ_u			2.438**	0.210	2.355**	0.205
Log Likelihood				-1311.265		-1281.752**

* Significant at .05-level; ** Significant at .01-level (2-sided).

[†]Hypothesized effect not significant when only the first treatment is considered.

significant, confirming Hypothesis 5. Also, subjects are more likely to reward the more the recipient contributes, but not the higher the positive deviation from the average. We do find that rewarding is less likely the more the recipient negatively deviates. Hypothesis 3 is supported: subjects who made a higher contribution are more likely to reward.

Model 9 shows the interaction of a subjects' own contribution with the contribution and positive deviation of the recipient. The significant effects indicate that high contributors are more likely to reward the higher and the further above the average someone contributes, confirming Hypotheses 4a and 4b.

The main effect of contribution (Hypothesis 3) and the effect that high contributors reward especially others who contribute much (Hypothesis 4a) are not found if we only consider the first treatment for Model 9. Again, this is probably due to the lower number of observations when only one treatment is included, which makes it more difficult to disentangle the different reasons why high contributors reward others.

5. Conclusion and discussion

We compared the effect of individual, majority, and unanimity decision rules for implementing punishment and reward on actors' ability to enforce cooperation in a Public Goods Game (PGG). For punishment, we conjectured that contributions are higher under a majority than an individual decision rule (Hypothesis 10a). However, we find higher contributions under the individual decision rule instead. As expected, we do find that contribution is lower under a unanimity than an individual punishment decision rule (Hypothesis 10b). For reward, the hypotheses concerning the effects of decision rules on contribution are all confirmed. We find that contribution is higher under an individual than a majority decision rule (Hypothesis 12a) and higher under a majority than a unanimity decision rule (Hypothesis 12b). In sum, for both punishment and reward contributions are lower, the more actors are required to agree on sanctioning. Also, as hypothesized, contribution is higher under punishment than reward for every decision rule (Hypothesis 11), although no difference is found in the majority condition when only the first treatment with sanctions is considered.

Findings on individual behavior, as captured in micro-level hypotheses, offer an explanation for the observed differences in contribution between decision rules. The emerging pattern is very similar for reward and punishment. Hypotheses on the use of cooperation-enforcing sanctions are all confirmed. High contributors are more likely to punish (Hypothesis 1) and to reward (Hypothesis 3) than low contributors. These high contributors enforce the norm that others should contribute as well. That is, they are more likely to punish the less a recipient contributes (Hypotheses 2a) and the lower the contribution of the recipient is compared to the other group members (Hypothesis 2b). Likewise, high contributors reward group members who also make a high contribution (Hypothesis 4a) and who contribute more compared to the others (Hypothesis 4b). In other words, there is

more consensus on sanctions among high contributors, the more an actor violates or adheres to their cooperative norm. Still, many punishments and rewards under the majority and unanimity decision rules were not executed. This implies that reaching the required number of actors was difficult despite the high consensus on whom to target.

When low contributors are punished, they contribute more in the subsequent interaction (Hypothesis 6). Similarly, actors who are rewarded for contributing more than other group members increase their contribution compared to others who are rewarded less (Hypothesis 8). Thus, we find strong evidence that cooperation-enforcing sanctions have a positive effect on contributions. Conversely, perverse sanctioning occurred too infrequently to affect contribution levels. We cannot confirm that high contributors decrease their contribution after being punished perversely (Hypothesis 7). Likewise, contrary to our expectations, free-riders who are rewarded perversely do not decrease their contribution further (Hypothesis 9). We did find that almost all perverse sanctions were ruled out under majority and unanimity.

In sum, we find strong evidence for cooperation-enforcing sanctions, and their positive effects on contribution. Concurrently, perverse sanctions occur too infrequently to affect cooperation. This makes an individual decision rule (IDR) unproblematic: punishment is mostly targeted at free-riders regardless of the possibility for individual actors to sanction perversely. Because more cooperation-enforcing sanctions are obstructed the more actors are required for the collective decision rule (CDRs), we observe lower contribution levels the more actors are required to agree. The observed micro-level behavior thus explains the macro-level finding of lower contribution levels under unanimity than majority, and lower contributions in the majority than in the individual condition.

The use of rewards decreases over time (Hypothesis 5). This provides an additional impediment for CDRs, because it implies that the more actors are required to agree, the sooner consensus cannot be reached anymore. Rewards are therefore even more problematic to enforce than punishment, hence contributions are higher under punishment than reward.

Casari and Luini (2009) find, in groups of five, that punishments on which two out of four actors agree are much more effective than sanctions with an IDR. We use stricter CDRs of two and three actors in groups of four, and find that contributions are highest under an IDR. However, contribution levels in our majority punishment condition (Figure 1) and the CDR of Casari and Luini (2009, Figure 1) are very similar. The difference between their findings and ours is that contribution in their individual punishment condition is much lower. Herrmann et al. (2008) find such differences in contributions under individual decision rules between subject pools. They attribute this to different levels of perverse punishment. Indeed, Casari and Luini (2009) find that contributions in their individual punishment condition are diminished due to perverse punishments. We find that perverse punishments do not affect contributions even under an IDR.

We started this paper with the observation that actors engaged in real-life public good problems often use CDRs to successfully enforce cooperation. One possible reason why we find that an IDR is more effective might be that interactions in our experiment are one-shot and anonymous rather than repeated. In many real-life public good problems, especially in small communities or between nations, participants interact repeatedly. Moreover, actors can often communicate before deciding whether or not to sanction. Repeated interaction and communication both imply that actors can coordinate on raising the required proportion of agreeing actors.

Furthermore, in real-life it is often possible to identify which actors neglected to agree on sanctioning. Therefore, when the required consensus is not reached the actors who did not sanction can be held accountable, for example through second-order punishment (Cinyabuguma et al. 2006; Denant-Boemont et al. 2007; Nikiforakis 2008). Previous studies found that second-order punishment is not always effective because it is used by defectors to punish first-order punishers. This issue should be alleviated when CDRs are used, because responsibility for punishment is shared by multiple others. Also, when a CDR is used for second-order punishment as well, agreement on punishment of punishers might not be reached.

Finally, in our experiment actors had complete information about others' contributions. In reality, some of the actors might make an inaccurate observation of the contributions of some of the others. An IDR might lead to inaccurate sanctioning decisions in such an environment (Grechenig et al. 2010; Ambrus and Greiner 2012). However, under a CDR mistaken sanctions caused by a wrong observation of an actor's contribution by one of the others will be ruled out.

Repeated interactions, communication on whom to sanction, public announcement of sanctioning decisions, use of counter-punishment, and noise can be implemented in future experiments to enhance resemblance with actual public good problems. As indicated above, these adaptations might favor CDRs, because coordination of sanctions in CDRs can become easier and mistakes in sanctions can be prevented. Still, the disentangling of sanctions through reciprocal contributions and sanctions through exogenous institutions will remain a challenge in these set-ups. In addition, there might also be some more realistic specifications of the interaction situation which favor IDRs. Most importantly, we assumed that non-implemented sanctions are costless. In reality, it might be more plausible that people have to invest in sanctioning before knowing whether others will agree. This would make implementation of sanctions under a CDR even more problematic. Future research should further specify conditions under which either CDRs or IDRs are more successful in enforcing cooperation.

Literature cited

- Ambrus, A. and B. Greiner. 2012. Imperfect Public Monitoring with Costly Punishment – An Experimental Study. *American Economic Review* 102:3317–3332.

- Anderson, C. M. and L. Putterman. 2006. Do Non-strategic Sanctions Obey the Law of Demand? The Demand for Punishment in the Voluntary Contribution Mechanism. *Games and Economic Behavior* 54:1–24.
- Andreoni, J. and L. K. Gee. 2012. Gun for Hire: Delegated Enforcement and Peer Punishment in Public Goods Provision. *Journal of Public Economics* 96:1036–1046.
- Balliet, D., L. B. Mulder, and P. A. M. Van Lange. 2011. Reward, Punishment, and Cooperation: A Meta-analysis. *Psychological Bulletin* 137:594–615.
- Bochet, O., T. Page, and L. Putterman. 2006. Communication and Punishment in voluntary Contribution Experiments. *Journal of Economic Behavior and Organization* 60:11–26.
- Botelho, A., G. W. Harrison, L. M. C. Pinto, and E. E. Rutström. 2005. Social Norms and Social Choice. Working Paper.
- Botelho, A., G. W. Harrison, L. M. Pinto, and E. E. Rutström. 2009. Testing Static Game Theory with Dynamic Experiments: A Case Study of Public Goods. *Games and Economic Behavior* 67:253–265e.3.
- Buchanan, J. M. and G. Tullock. 1962. *The Calculus of Consent*. Ann Arbor: The University of Michigan Press.
- Carpenter, J. P. 2007. The Demand for Punishment. *Journal of Economic Behavior and Organization* 4:522–542.
- Carpenter, J. and P. H. Matthews. 2009. What Norms Trigger Punishment? *Experimental Economics* 12:272–288.
- Casari, M. and L. Luini. 2009. Cooperation Under Alternative Punishment Institutions: An experiment. *Journal of Economic Behavior and Organization* 71:273–282.
- Chaudhuri, A. 2011. Sustaining Cooperation in Laboratory Public Goods Experiments: A selective Survey of the Literature. *Experimental Economics* 14:47–83.
- Choi, J.-K. and T. K. Ahn. 2013. Strategic Reward and Altruistic Punishment Support Cooperation in a Public Goods Game Experiment. *Journal of Economic Psychology* 35:17–30.
- Cinyabuguma, M., T. Page, and L. Putterman. 2006. Can Second-order Punishment Deter Perverse Punishment? *Experimental Economics* 9:265–279.
- Dari-Mattiacci, G. and G. Geest. 2010. Carrots, Sticks, and the Multiplication Effect. *Journal of Law, Economics and Organization* 26:365–384.
- Dawes, R. M. 1980. Social Dilemmas. *Annual Review of Psychology* 31:169–193.
- Dawes, C. T., J. H. Fowler, T. Johnson, R. McElreath, and O. Smirnov. 2007. Egalitarian Motives in Humans. *Nature* 446:794–796.
- Decker, T., A. Stiehler, and M. Strobel. 2003. A Comparison of Punishment Rules in Repeated Public Goods Games – An Experimental Study. *Journal of Conflict Resolution* 47:751–772.
- Denant-Boemont, L., D. Masclet, and C. N. Noussair. 2007. Punishment, Counterpunishment and Sanction Enforcement in a Social Dilemma Experiment. *Economic Theory* 33:145–167.

- De Quervain, D. J. F., U. Fischbacher, V. Treyer, M. Schellhammer, U. Schnyder, A. Buck, and E. Fehr. 2004. The Neural Basis for Altruistic Punishment. *Science* 305:1254–1258.
- Diekmann, A. and T. Voss. 2003. Social Norms and Reciprocity. Discussion paper, University of Leipzig.
- Drouvelis, M. and J. C. Jamison. 2012. Selecting Public Goods Institutions: Who Likes to Punish and Reward? Working paper, 60.
- Ellingsen, T., B. Herrmann, M. A. Nowak, D. G. Rand, and C. E. Tarnita 2012. Civic Capital in Two Cultures: The Nature of Cooperation in Romania and USA. SSRN working paper, 60.
- Ertan, A., T. Page, and L. Putterman. 2009. Who to Punish? Individual Decisions and Majority Rule in Mitigating the Free Rider Problem. *European Economic Review* 38:495–511.
- Fehr, E. and U. Fischbacher. 2004. Social Norms and Human Cooperation. *TRENDS in Cognitive Sciences* 8:185–190.
- Fehr, E. and S. Gächter. 2000. Cooperation and Punishment in Public Goods Experiments. *American Economic Review* 90:980–994.
- Fehr, E. and S. Gächter. 2002. Altruistic Punishment in Humans. *Nature* 415:137–140.
- Fehr, E. and H. Gintis. 2007. Human Motivation and Social Cooperation: Analytical and Experimental Foundations. *Annual Review of Sociology* 33:43–64.
- Fischbacher, U. 2007. z-Tree – Zurich Toolbox for Readymade Economic Experiments – Experimenter’s Manual. *Experimental Economics* 10:171–178.
- Fischbacher, U. and S. Gächter. 2010. Social Preferences, Beliefs, and the Dynamics of Free Riding in Public Good Experiments. *American Economic Review* 100:541–556.
- Fischbacher, U., S. Gächter, and E. Fehr. 2001. Are People Conditionally Cooperative? Evidence from a Public Goods Experiment. *Economics Letters* 71:397–404.
- Fudenberg, D. and P. A. Pathak. 2010. Unobserved Punishment Supports Cooperation. *Journal of Public Economics* 94:78–86.
- Gächter, S. 2007. Conditional Cooperation. Behavioral Regularities from the Lab and the Field and their Policy Implications. In *Economics and Psychology. A Promising New Cross-disciplinary Field*, eds. B. Frey and A. Stutzer. Cambridge: MIT Press.
- Gächter, S. and B. Herrmann. 2009. Reciprocity, Culture and Human Cooperation: Previous Insights and a New Cross-cultural Experiment. *Philosophical Transactions of the Royal Society* 364:791–806.
- Gächter, S. and C. Thöni. 2011. Micromotives, Microstructure and Macrobehavior: The Case of Voluntary Cooperation. *Journal of Mathematical Sociology* 35:26–65.
- Grechenig, C., A. Niklisch, and C. Thöni. 2010. Punishment Despite Reasonable Doubt – A Public Goods Experiment with Sanctions Under Uncertainty. *Journal of Empirical Legal Studies* 7:847–867.

- Greiner, B. 2004. *The Online Recruitment System ORSEE 2.0. A guide for the organization of experiments in economics*. University of Cologne, Working Paper Series in Economics 10.
- Guillen, P., C. Schwieren, and G. Staffiero. 2006. Why Feed the Leviathan? *Public Choice* 130:115–128.
- Herrmann, B. and C. Thöni. 2009. Measuring Conditional Cooperation: A Replication Study in Russia. *Experimental Economics* 12:87–92.
- Herrmann, B., C. Thöni, and S. Gächter. 2008. Perverse Punishment Across Societies. *Science* 319:1362–1367.
- Isaac, M. R. and J. M. Walker. 1988. Communication and Free-riding Behavior: The Voluntary Contribution Mechanism. *Economic Inquiry* 26:585–608.
- Kamei, K. 2012. From Locality to Continent: A Comment on the Generalization of an Experimental Study. *The Journal of Socio-Economics* 41:207–210.
- Kamei, K., L. Putterman, and J.-R. Tyran. 2011. State or Nature? Formal vs. Informal Sanctioning in the Voluntary Provision of Public Goods. SSRN working Paper, 49.
- Kocher, M. G., P. Martinsson, and M. Visser. 2007. Does Stake Size Matter for Cooperation and Punishment? *Economic Letters* 99:508–511.
- Kocher, M. G., T. Cherry, S. Kroll, R. J. Netzer, and M. Sutter. 2008. Conditional Cooperation on Three Continents. *Economics Letters* 101:175–178.
- Ledyard, J. O. 1995. Public Goods: A Survey of Experimental Research. In: *Handbook of Experimental Economics*, eds. J. Kagel and A. E. Roth, 111–194. New York: Princeton University Press.
- Markussen, T., L. Putterman, and J.-R. Tyran. 2011. Self-organization for Collective Action: An Experimental Study of Voting on Formal, Informal, and No Sanction Regimes. SSRN working paper, 49.
- Masclot, D., C. Noussair, S. Tucker, and M.-C. Villeval. 2003. Monetary and Nonmonetary Punishment in the Voluntary Contributions Mechanism. *The American Economic Review* 93:366–380.
- Milinski, M. and B. Rockenbach. 2011. On the Interaction Between the Stick and the Carrot in Social Dilemmas. *Journal of Theoretical Biology* 229:139–143.
- Monin, B. 2007. Holier than Me? Threatening Social Comparison in the Moral Domain. *International Review of Social Psychology* 20:53–68.
- Nikiforakis, N. 2008. Punishment and Counter-punishment in Public Good Games: Can we Really Govern Ourselves? *Journal of Public Economics* 92:91–112.
- Nikiforakis, N. and H.-T. Normann. 2008. A Comparative Statics Analysis of Punishment in Public-good Experiments. *Experimental Economics* 11:358–369.
- Nosenzo, D. and M. Sefton. 2012. Promoting Cooperation: The Distribution of Reward and Punishment Power. *CeDEx Discussion Paper Series*, 32.
- Olson, M. 1965. *The Logic of Collective Action. Public Goods and the Theory of Groups*. Cambridge: Harvard University Press.
- Ones, U. and L. Putterman. 2007. The Ecology of Collective Action: A Public Goods and Sanctions Experiment with Controlled Group Formation. *Journal of Economic Behavior & Organization* 62:495–521.

- Ostrom, E. 1990. *Governing the Commons*. Cambridge: Cambridge University Press.
- Ostrom, E. 2000. Collective Action and the Evolution of Social Norms. *The Journal of Economic Perspectives* 14:137–158.
- Ostrom, E. 2010. Beyond Markets and States: Polycentric Governance of Complex Economic Systems. *American Economic Review* 100:1–33.
- Ostrom, E. 2012. Coevolving Relationships Between Political Science and Economics. Mimeo, Workshop on Political Theory and Policy Analysis, Indiana University.
- Ostrom, E., J. Walker, and R. Gardner. 1992. Covenants with and without a Sword: Self-governance is Possible. *The American Political Science Review* 86:404–417.
- Putterman, L., J.-R. Tyran, and K. Kamei. 2011. Public Goods and Voting on Formal Sanction Schemes. *Journal of Public Economics* 95:1213–1222.
- Rand, D. G. and M. A. Nowak. 2011. The Evolution of Antisocial Punishment in Optional Public Goods Games. *Nature Communications* 2:434.
- Rand, D. G., A. Dreber, T. Ellingsen, D. Fudenberg, and M. A. Nowak. 2009. Positive Interactions Promote Public Cooperation. *Science* 325:1272–1275.
- Rand, D. G., J. J. Armao, M. Nakamaru, and H. Ohtsuki. 2010. Anti-social Punishment Can Prevent the Co-evolution of Punishment and Cooperation. *Journal of Theoretical Biology* 265:624–632.
- Sefton, M., R. Shupp, and J. M. Walker. 2007. The Effects of Rewards and Sanctions in Provision of Public Goods. *Economic Inquiry* 45:671–690.
- Selten, R. 1967. Die Strategiemethode zur Erforschung des eingeschränkt rationalen Verhaltens im Rahmen eines Oligopolexperimentes. In: *Beiträge Zur Experimentellen Wirtschaftsforschung*, eds. H. Saueremann, 136–168. Tübingen: J.C.B. Mohr.
- Sutter, M., S. Haigner, and M. G. Kocher. 2010. Choosing the Carrot or the Stick? Endogenous Institutional Choice in Social Dilemma Situations. *Review of Economic Studies* 77:1540–1566.
- Thöni, C., J.-R. Tyran, and E. Wengström. 2012. Microfoundations of Social Capital. *Journal of Public Economics* 96:635–643.
- Traulsen, A., T. Röhl, and M. Milinski. 2012. An Economic Experiment Reveals that Humans Prefer Pool Punishment to Maintain the Commons. *Proceedings of the Royal Society B* 279:3716–3721.
- Tyran, J.-R. and L. P. Feld. 2006. Achieving Compliance when Legal Sanctions are Non-deterrent. *Scandinavian Journal of Economics* 108:135–156.
- Van Lange, P. A. M. 1999. The Pursuit of Joint Outcomes and Equality in Outcomes: An Integrative Model of Social Value Orientation. *Journal of Personality and Social Psychology* 77:337–349.
- Veszteg, R. F. and E. Narhetali. 2010. Public-good Games and the Balinese. *International Journal of Social Economics* 37:660–675.
- Vyrastekova, J. and D. Van Soest. 2008. On the (In)effectiveness of Rewards in Sustaining Cooperation. *Experimental Economics* 11:53–65.

- Walker, J. M. and M. A. Halloran. 2004. Rewards and Sanctions and the Provision of Public Goods in One-Shot Settings. *Experimental Economics* 7:235–247.
- Wiedemann, V., D. Barrera, and V. Buskens. 2011. The Consequences of Monetary Rewards and Punishment on Cooperation in Repeated Public Goods Games. SSRN working paper, 19.
- www.europa.eu, accessed at 03-11-2010.
- www.un.org, accessed at 03-11-2010.
- Yamagishi, T. 1986. The Provision of a Sanctioning System as a Public Good. *Journal of Personality and Social Psychology* 51:110–116.