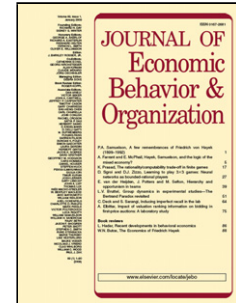


## Accepted Manuscript

Title: An individualistic approach to institution formation in public good games

Author: Abhijit Ramalingam Sara Godoy Antonio J. Morales  
James M. Walker



PII: S0167-2681(16)30111-1  
DOI: <http://dx.doi.org/doi:10.1016/j.jebo.2016.06.003>  
Reference: JEBO 3821

To appear in: *Journal of Economic Behavior & Organization*

Received date: 13-6-2015  
Revised date: 26-5-2016  
Accepted date: 6-6-2016

Please cite this article as: Ramalingam, Abhijit, Godoy, Sara, Morales, Antonio J., Walker, James M., An individualistic approach to institution formation in public good games. *Journal of Economic Behavior and Organization* <http://dx.doi.org/10.1016/j.jebo.2016.06.003>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

# An individualistic approach to institution formation in public good games

Abhijit Ramalingam <sup>a\*</sup>, Sara Godoy <sup>b</sup>, Antonio J. Morales <sup>c</sup>, James M. Walker <sup>d</sup>

<sup>a</sup> School of Economics and Centre for Behavioural and Experimental Social Science, University of East Anglia, Norwich NR4 7TJ, UK, a.ramalingam@uea.ac.uk

<sup>b</sup> Centre for Behavioural and Experimental Social Science, University of East Anglia, Norwich NR4 7TJ, UK, s.godoy-garzon@uea.ac.uk

<sup>c</sup> Facultad de Economía, Universidad de Málaga, Málaga 29007, Spain, amorales@uma.es

<sup>d</sup> Department of Economics and Vincent and Elinor Ostrom Workshop in Political Theory and Policy Analysis, Indiana University, Bloomington IN 47405, USA, walkerj@indiana.edu

\*Corresponding author: abhi.ramalingam@gmail.com, Tel: +44-1603-597382, Fax: +44-1603-456259.

May 26, 2016

**Highlights:**

- In a public goods experiment, we investigate if individuals will unilaterally provide a sanctioning institution
- Only individuals who give themselves the ‘right’ to punish can do so
- Sanctioning need not be provided at group level; individuals can, and do, provide sanctioning in groups
- Sanctioning is effective at raising cooperation when its provision is costless
- However, even a negligible monetary provision cost leads to counterproductive use of sanctioning

**Abstract**

In a repeated public goods setting, we explore whether individuals, acting unilaterally, will provide an effective sanctioning institution. Subjects first choose independently whether they will participate in a sanctioning stage that follows a contribution stage. Only those who gave themselves the “right” to sanction can do so. We find that the effectiveness of the institution may not require provision of the institution at the level of the group. Individuals acting unilaterally are able to provide sanctioning institutions that effectively raise cooperation. The effectiveness of the institution, however, depends on whether the “right” to sanction entails a monetary cost or not.

**Keywords:** public goods, experiment, punishment, institution formation, unilateral provision, cooperation

**JEL codes:** C72, C91, C92, D02, H41

## 1. Introduction

A common feature of most studies of public good games is that the institution proposed to increase contributions is provided exogenously and the emphasis is placed on the conditions which effectively help to alleviate the free rider problem (see Chaudhuri, 2011 for a recent survey). Of late, there is growing interest in how the institution comes into being. This issue is important because the formation of the institution is subject to a second-order free rider problem. Others may profit from the institution but they prefer someone else to provide it (see Oliver, 1980).<sup>1</sup> The literature on the endogenous formation of institutions provides an answer assuming that the institutional choice mechanism is voting: there is ample experimental evidence showing that in many cases, the outcome of voting is a sanctioning institution.<sup>2</sup> This approach however assumes that the group has the capacity to organize the voting mechanism and to enforce the resulting outcome.

Individuals in many societies can and do act on their own – such as deciding on contributions to the public good – without the need for the group to aggregate individual preferences. In addition, in many settings, individuals discontented with the contribution levels of their peers, can choose to unilaterally provide and enforce sanctioning institutions.<sup>3</sup> It is, therefore, perhaps more natural to take individual actions as the starting point in analysing the ability of groups to endogenously provide and enforce potentially efficiency-enhancing institutions such as sanctioning.

In this paper, we experimentally examine the provision and effectiveness of the sanctioning institution in a public goods game when its provision depends on *individuals* acting independently. Will individuals unilaterally choose a punishment role? If so, what is the effect on group outcomes in comparison to when the sanctioning institution is exogenously and

---

<sup>1</sup> Other early works on this issue are Yamagashi (1986), where subjects were offered the possibility to voluntarily fund a sanctioning institution in a public goods game and Ostrom et al. (1992), where in a common pool resource game, subjects had the opportunity to communicate to decide whether to use sanctions. Traulsen et al. (2012) and Zhang et al. (2014) find that many players choose pool punishment when second order punishment of non-punishers is possible.

<sup>2</sup> Gürer et al. (2006), Ertan et al. (2009) and Sutter et al. (2010) are examples where the choice is between no sanctions versus informal sanctions. In Kosfeld et al. (2009) and Kube et al. (2015), the choice is between no sanctions and formal sanctions imposed by a central authority. Markussen et al. (2014) and Kamei et al. (2015) are recent studies where the choice is between formal and informal sanction schemes.

<sup>3</sup> There are alternative institutions other than sanctioning that can be implemented. Some examples are rewards for high contributors (Sefton et al. 2007), ostracism of low contributors (Cinyabuguma et al. 2005), excludability (Croson et al., 2014), leadership within groups (van der Heijden et al. 2009) and formation of coalitions (Dannenberg et al., 2010 and McEvoy et al., 2011). Kube et al. (2015) study the endogenous provision of institutions that include both minimum contribution levels and centralised sanctioning.

universally provided? Finally, how is the effectiveness of the institution changed if individuals must unilaterally bear the cost of providing it?

In our experiment, before making decisions on contributions, individuals unilaterally decide whether or not they want to be able to use punishment. The number of such individuals is then announced before the contribution stage takes place. Finally, contribution levels are made public and only those individuals who gave themselves the “right” to make use of sanctioning can assign punishment to *any* group member.<sup>4</sup> This is akin to the behaviour of vigilantes who take it upon themselves to provide mechanisms to enforce a norm and punish others who violate it or some voluntary neighbourhood watch groups that provide both monitoring and sanctioning. An obvious behavioural question is whether group members will respond differently to a sanctioning mechanism that has been exogenously provided to all group members in comparison to one in which individuals act unilaterally to choose to provide the mechanism.

We consider two variants of the sanctioning institution where individuals choose-to-participate (CTP) - whether the choice to participate is available at no monetary cost (CTP0) or whether there is a positive cost (CTP1).<sup>5</sup> In addition, we replicate the most common settings in public goods experiments – the Voluntary Contributions Mechanism (VCM) and the VCM with an exogenously provided opportunity to punish (StdPun). In the VCM setting, subjects could only contribute to the public good and there was no enforcement mechanism available. In StdPun, all group members automatically had the right to assign punishment to others in the group.

Based on the standard assumption of own income maximization, individuals would not be expected to provide the sanctioning institution or to use it to discipline free-riders. However, previous work has found that individuals do make use of exogenously provided sanctioning institutions and are able to enforce high cooperation levels in groups. Fehr and Schmidt (1999), hereafter FS, rationalise such behaviour using a model of inequity aversion. Extending their model to our setting, we find that, as in FS, any symmetric contribution profile can be supported

---

<sup>4</sup> A related paper is Masclet et al. (2013), where subjects can make non-binding threats before the contribution stage. Players issue costless detailed threats to other group members as a function of hypothetical contribution levels and these threats are made public before making contribution decisions. They find an increase in contributions relative to a standard VCM.

<sup>5</sup> Using standard economic terminology, the punishment technology may entail a *fixed* per round provision cost associated with acquiring and having the technology ready to use, and a *variable* cost associated with making use of it. The standard approach in the literature is linear variable cost with no provision cost (as in Herrmann et al., 2008). Some papers, though, consider a positive provision cost but the decision to provide the sanctioning institution is taken at the group level (see for example Kosfeld et al., 2009).

as a subgame perfect Nash equilibrium. However, this requires everyone in the group – selfish and inequity averse players – to provide the sanctioning institution. In addition, there exist subgame perfect equilibria with less than complete provision where only a subset of inequity averse players provide the institution. However, to account for the pecuniary inequity that arises from the different participation decisions, contribution profiles in such equilibria are asymmetric.

One may think of the CTP settings as allowing for *extreme* cases that correspond to the provision cost of the sanctioning institution. When the provision cost approaches infinity, no player will choose to sanction and the institution will resemble the VCM. When the provision cost approaches zero as in CTP0, then all players may choose to give themselves the right to sanction and the institution will resemble StdPun. In the intermediate range however, the FS model predicts a multiplicity of equilibrium outcomes with a wide range of contribution levels and participation in the punishment stage. It is such situations that our experiment allows us to investigate.

Our experimental data shows several monotonic results. When the provision of the sanctioning institution is costly, fewer subjects choose to participate in the punishment stage than when it is costless. In terms of the effects on cooperation, while both CTP treatments start at the same level, cooperation levels in the two CTP treatments soon diverge. In CTP0, groups are as successful in raising cooperation as with automatic universal participation in punishment (StdPun). In CTP1, despite the fact that some players do provide the sanctioning institution, groups are unable to raise cooperation levels and contributions to the public good stagnate at levels close to those observed in the VCM setting. However, complementary to these general patterns, there are a number of additional findings that greatly enrich the picture.

First, in the costless treatment, there is less than full provision of the sanctioning institution; in only 10% of all occasions did all group members choose to participate in the punishment stage and the overall average participation rate is 60%. The literature on voting on punishment systems in public good games sheds some light on this result (see for example, Gurerk et al., 2006, and Ertan et al., 2009). First, not all groups succeed in implementing the punishment regime and second, in those cases that the group implements the punishment system, the institution is not always unanimously approved (majority rule is usually used).<sup>6</sup> This means

---

<sup>6</sup> Gerber et al. (2013) find that a unanimous participation rule is the most effective at increasing the number of implemented institutions.

that some subjects are not in favour of sanctioning others as an institution. In our CTP settings, these subjects (selfish players that will never find it rational to punish) may choose to not take on the punisher role, leaving that role for the inequity averse players that will ultimately be responsible for making the threats to punish credible.<sup>7</sup>

Second, monotonicity holds between *and* within the CTP settings. On one hand, the average number of subjects providing the institution is larger in CTP0 than in CTP1. This suggests that the law of demand previously reported in the literature with respect to variable punishment costs (Anderson and Putterman, 2006, and Carpenter, 2007) extends to fixed provision cost.<sup>8</sup> On the other hand, within each CTP treatment, there is a positive relation between the number of players choosing to provide the institution and group contribution levels. This result suggests that the threat to punish was credible in both CTP treatments.

Third, controlling for the number of participants, contributions in CTP0 are higher than in CTP1. The question is why the development of credible punishment threats increases contributions in the costless setting to a larger extent than in the costly one. In CTP1, the participation decision is strongly contingent on having been punished in the previous round. This is not true in CTP0. Further, in regard to the use of sanctioning, subjects are found to punish high and low contributors with virtually the same intensity in CTP1, but not in StdPun or in CTP0. This suggests that “blind revenge” (Ostrom et al. 1992) is a larger factor in CTP1, diminishing the efficacy of targeted punishment of low contributors, the key element for raising contributions.

Fourth, an individual’s decision to provide the sanctioning institution is not found to be strongly correlated with his/her contribution decision. This suggests that individuals’ cooperation decisions depend more on the persistent existence of a sanctioning institution and less so on whether they themselves provide the institution.

---

<sup>7</sup> There is a branch of the literature that analyses the performance of the sanctioning institution in VCM settings where *exogenously provided* punishment networks limit punishment opportunities, as well as the information subjects receive on contributions and punishment imposed/received (e.g., Carpenter et al., 2012, Fatas et al., 2010 and Leibbrandt et al., 2015). Carpenter et al. (2012) find that the complete network, where everybody can punish everyone, is more efficient than incomplete networks that restrict punishment opportunities to a subset of subjects. Leibbrandt et al. (2015) examine complete vs. incomplete punishment networks, but in a setting where there are fixed identifiers across rounds that allow subjects to receive complete information about all other subjects in their group regarding contribution and punishment decisions. They find that the structure of the punishment network significantly affects allocations to the public good and that network configurations are more important than punishment capacities.

<sup>8</sup> Although in a different context, there are studies showing that zero is a special prize, in the sense that people perceive the benefits associated with free products as higher (Shampanier et al., 2007).

Finally, in CTP1, the experimental value of the cost of providing the sanctioning institution was negligible - a twentieth of an individual's initial endowment. After completing the initial experiments mentioned above, we conducted a variant of the costly treatment, CTP5, in which the cost of participation was higher than in CTP1. We find essentially the same patterns in punishment and cooperation behaviour as in CTP1. This suggests that the *mere existence* of a provision cost hinders the development of an *effective* sanctioning institution. The reason for this result appears to be related to both a decrease in the level of participation in, and use of, the sanctioning institution.

To our knowledge, no previous study explicitly examines treatment conditions with both positive and null provision costs of providing a sanctioning institution. There is, however, some prior indirect evidence. Both in Gürer et al. (2006) - where players can vote with their feet whether to be in a society with or without punishment - and in Ertan et al. (2009) - where the group decides whether punishment is allowed using a majority rule - the provision cost of the sanctioning institution is zero and it is effectively chosen with positive effects on contributions and efficiency levels. Kosfeld et al. (2009) consider a positive provision cost in a setting in which players voted for implementing the institution. The provision cost, however, is borne by only those who voted for provision. They find that punishment is successfully implemented by a large number of groups.<sup>9</sup>

Our results indicate that an endogenous sanctioning institution can raise contributions, even without full provision. The persistent participation of players (in CTP0, the average participation rates of the players with the first and second highest number of decisions to participate are 93% and 81%) and punishment targeted at low contributors are found to be behind the successful implementation of the institution. However, our results also suggest that endogenous institutional change can be a very fragile process that is sensitive to subtle institutional details; in our case, to the existence of a positive provision cost.

The rest of the paper is organised as follows. Section 2 details our experimental design and procedures. Section 3 theoretically explores the effects of a participation cost and presents our hypotheses. Section 4 presents and discusses our results, including results from the additional treatment designed to explore the effect of a more substantial participation cost, and Section 5

---

<sup>9</sup> In studies where the subjects' choice is between formal and informal sanctioning (Markussen et al., 2014, and Kamei et al., 2015), the cost of providing the formal mechanism affects the choice: formal sanctions are more popular when they carry no up-front cost, whereas informal sanctions are more popular and efficient when adopting the formal scheme entails such a cost.



concludes. Appendix A of the online material contains the experimental instructions and Appendix B contains additional analysis used to support results presented in the paper. Appendix C presents an analysis of individuals' final earnings in relation to the inequity-aversion model.

## 2. Experimental Design and Procedures

The base game in our experiment was a standard Voluntary Contributions Mechanism (VCM). In the VCM, a group of  $n \geq 2$  players repeatedly plays a linear public goods game with just one stage - a contribution stage. Each player  $i$  ( $i = 1, 2, \dots, n$ ) begins each round with an endowment of  $y$  tokens in a private account from which he/she can allocate  $g_i \in \{0, 1, 2, \dots, y\}$  to a group account, i.e., the public good. The balance,  $e_i = y - g_i$ , remains in the private account and earns a return of 1. Each player in the group receives  $aG$  from the group account where  $G = \sum_{i=1}^n g_i$  is the total contribution to the public good and  $a$  ( $0 < a < 1 < an$ ). The monetary payoffs to player  $i$  in a round are given by

$$\pi_i(\mathbf{g}) = (y - g_i) + aG$$

where  $\mathbf{g}$  is the profile of contributions to the public good. The Nash equilibrium in the stage game is for each player to contribute nothing to the public good ( $g_i = 0 \forall i = 1, 2, \dots, n$ ) while the social optimum is for each player to contribute his/her entire endowment to the public good ( $g_i^* = y \forall i = 1, 2, \dots, n$ ). The Nash equilibrium and the social optimum remain unchanged under finite repetitions of the stage game.

In games with punishment, players can also use their earnings from the contribution stage of the game to punish other players in a subsequent punishment stage. Let  $p_{kl}$  denote the punishment player  $k$  sends to player  $l$ ,  $k \neq l$ . A unit of punishment imposed on a player costs the punishing player  $c$  units ( $0 < c < 1$ ).<sup>10</sup> Denoting the punishment profile in the group by  $\mathbf{p}$ , a player's monetary payoff in a period is given by

$$\pi_i(\mathbf{g}, \mathbf{p}) = (y - g_i) + aG - c \sum_{\substack{j=1 \\ j \neq i}}^n p_{ij} - \sum_{\substack{j=1 \\ j \neq i}}^n p_{ji}.$$

---

<sup>10</sup> We use the notation from FS, in particular regarding the description of the punishment technology.

Individual contributions to the public good at the Nash equilibrium and at the social optimum in the punishment game (in one-shot and finitely repeated games) are identical to those in the game without punishment, i.e., zero and full contributions respectively. In addition,  $p_{kl} = 0 \forall k, l$  at both the Nash equilibrium and at the social optimum.

In all treatments, there were 20 rounds with fixed groups and a contribution stage with  $n = 4$ ,  $y = 20$ , and  $a = 0.5$ . At the end of the contribution stage, each subject was informed of her group's total contribution to the public good in that round, the individual contributions of the others in her group in descending order and her individual earnings from her private account and from the public good. Subjects did not have individual identifiers that could create reputation effects.

In the first treatment (VCM), a round ended after the contribution stage. In all other treatments, subjects played a punishment game after the contribution stage. The second treatment was the standard exogenously provided sanctioning institution (StdPun), as in Gächter et al. (2008). In this treatment, after the contribution stage, subjects could use their earnings from the contribution stage to reduce the earnings of each other, up to a maximum of 5 tokens for each other group member.<sup>11</sup> The term punishment was not used. For brevity here, however, we will refer to such reductions as punishment. All four subjects in a group automatically entered this stage, where they decided how much punishment to assign, if any, to each of the others in their group. Thus, while the assignment of punishment was endogenous, participation in the institution itself was exogenously imposed for *all* group members, and at no cost. The punishment technology used was 1:3, i.e., one token used to punish a group member cost the punishing member 1 token and the recipient 3 tokens (i.e.,  $c = 1/3$  in terms of FS notation). The costs of assigning and receiving punishment were deducted from earnings from the contribution stage.<sup>12</sup> After the punishment stage, subjects were informed of the *total* amount of punishment they received and their earnings from both stages of the round. Because no subject identifiers were used, subjects could not associate punishment received with the particular group member who assigned the punishment.

The CTP treatments endogenised the provision of the sanctioning institution. Each group member was required to choose, in each round, whether or not to provide the sanctioning

---

<sup>11</sup> cf. Sefton et al. (2007) where subjects were given an additional endowment for punishment.

<sup>12</sup> If a player's earnings from the contribution stage was lower than 15 tokens, punishment was limited by his earnings. A player could have negative earnings in a round, but could not earn negative amounts in the experiment.

institution, i.e. to participate in the punishment stage in a round. Prior to the contribution stage, each subject chose whether to participate in the punishment stage that followed the contribution stage.<sup>13</sup> Subjects had to pay a fee,  $\gamma \geq 0$ , to provide the institution. Before making contribution decisions, subjects were informed only of the number of people in their group who had chosen to participate in the punishment stage. Only those who indicated a willingness to participate in the punishment stage in a round could assign punishment after the contribution stage in that round. These subjects could then punish *any* other group member, i.e., all group members could receive punishment, regardless of their choice in the initial stage. If no subject in a group chose to participate in the punishment stage in a round, the round ended after the contribution stage.

The study included two initial CTP treatments. In CTP0, the decision to participate in the punishment stage was costless ( $\gamma = 0$ ) and the institution was provided for free to each group member who chose to participate. In CTP1, each group member choosing to participate in the punishment stage paid a fee of 1 token, i.e.,  $\gamma = 1$ . The fee was deducted from the earnings of the subject after the contribution stage and before the punishment stage. This was done to ensure that a subject who gave herself the right to punish could contribute as much to the public good as could a subject who chose not to participate in the punishment stage. The punishment technology-parameters were the same as in StdPun.

Table 1 summarises the initial treatments and presents the number of observations in each.

Fehr and Gächter (2000) show that punishment is significantly more effective in raising contributions to the public good under partner matching than under stranger matching. Hence, partner matching allows more room for exploration of the effects of intermediate treatments between VCM and StdPun. Further, each group forms an independent observation under partner matching while each session forms an independent observation under stranger matching. While avoiding the risk of potential contamination of subjects in an entire session, partner matching also allows us to maximise the number of independent observations.

All sessions were conducted at EssexLab at the University of Essex. In each session, 12 to 24 subjects, recruited from the student body at Essex were randomly and anonymously assigned to four-person groups that stayed fixed throughout the 20 rounds. The repeated nature of the game and the partner matching within groups was common information for all subjects. At the

---

<sup>13</sup> We used neutral language in the instructions and never referred to “contributions” or “punishment”. In Stage 1, subjects were asked “Do you want to make decisions in Stage 3?”

beginning of each session, instructions for the 20-round public goods game were read out by an experimenter. Subjects also had a copy of the instructions that they could refer to at any time during the experiment. Subjects then took a quiz to ensure understanding. They could not proceed until all questions were answered correctly. Subjects then made decisions privately at their computer terminals. At the end of the session, subjects answered a demographic questionnaire.

The experiment was programmed in z-Tree (Fischbacher, 2007). In all treatments, the stage game was repeated for 20 rounds and earnings from a round could not be carried forward to future rounds. Subjects were paid their earnings from all 20 rounds of the public goods game. Tokens were converted to Pounds at the rate of 60 tokens to £1. A session lasted about 55 minutes and subjects earned an average of £12.35 each including a £2.50 show-up fee.

### 3. Cooperation with Endogenous Institution Provision: Theoretical Predictions

Past studies have shown that cooperation can be sustained when the contribution stage is followed by a punishment stage. One approach to rationalising this finding is by using social preferences à la FS, that are defined in terms of final monetary outcomes. To guide in the analysis of the observed behaviour, we extend the FS setting to our endogenous provision game and seek to identify behavioural regularities that can generate testable hypotheses.

#### 3.1 Institution Provision and Punishment by Inequity Averse Individuals

FS consider groups composed of selfish and inequity averse players. For a profile of monetary payoffs  $(\pi_1, \dots, \pi_n)$ , the utility to an inequity-averse player  $i$  is

$$u_i(\pi_1, \dots, \pi_n) = \pi_i - \frac{\alpha_i}{n-1} \sum_{\substack{j=1 \\ j \neq i}}^n \max\{\pi_j - \pi_i, 0\} - \frac{\beta_i}{n-1} \sum_{\substack{j=1 \\ j \neq i}}^n \max\{\pi_i - \pi_j, 0\}$$

where  $\alpha_i$  measures the utility loss to player  $i$  associated with disadvantageous inequality and  $\beta_i$  measures the utility loss associated with advantageous inequality, with  $\alpha_i > \beta_i$  and  $\beta_i \in [0,1]$ . FS show that any symmetric contribution level  $g$  (ranging from 0 to full contribution) can be supported as a subgame perfect equilibrium outcome if there is a group of  $n'$  “conditionally cooperative enforcers”, with  $1 \leq n' \leq n$  such that:

(a)  $\beta_i \geq 1 - \alpha$ , and

(b)  $c < \frac{\alpha_i}{(n-1)(1+\alpha_i)-(n'-1)(\alpha_i+\beta_i)}$  for  $i = 1, 2, \dots, n'$ .

Condition (a) assures that enforcers dislike advantageous inequality and do not find contribution of zero to be a dominant strategy (they are conditionally cooperative). Condition (b) assures that enforcers experience a utility loss associated with disadvantageous inequality and thus find it optimal to punish free riders contributing  $g' < g$  by choosing  $p^* = \frac{g-g'}{n'-c}$ . This way, potential non-contributors find it optimal to contribute in the first stage.

In this section we extend the FS setting by adding an initial acquisition stage. The game is now composed of three stages:

Stage 1 - Acquisition decisions: players choose whether to acquire the punishment technology to be used in Stage 3. The acquisition cost is  $\gamma \geq 0$ .

Stage 2 - Contribution decisions: Upon observing decisions in Stage 1, players choose contributions levels.

Stage 3 - Punishment decisions: Upon observing decisions in Stage 2, those players who acquired the punishment technology in Stage 1 can punish *any* member of the group.

Compared to the FS setting, the sequential addition of a costly provision of punishment modifies two things. First, players can condition their contribution and punishment decisions on the participation decisions in Stage 1. We will use this feature to construct subgame perfect Nash equilibria that use the contribution stage to “punish” deviations from the prescribed participation decisions. Second, if a player has spent  $\gamma$  on acquiring the punishment technology, the monetary payoff (before punishment) to this player decreases by the amount  $\gamma$ . Given that inequity averse players base their punishment behaviour on relative concerns, asymmetries in participation decisions will carry over to the punishment stage when  $\gamma > 0$ , prompting punishment even if contributions are equal.

This second point above is crucial. Note that if player  $i$  has invested in the punishment technology and player  $j$  has not, then prior to any punishment, the monetary payoffs to players  $i$  and  $j$  are  $\pi_i(g) = (y - g_i - \gamma) + aG$  and  $\pi_j(g) = (y - g_j) + aG$ , respectively. And the payoff difference  $\pi_j - \pi_i$  amounts to  $(g_i + \gamma) - g_j$ . From the point of view of punishment behaviour, having spent  $\gamma$  on the acquisition of the punishment technology and contributed  $g_i$

is *equivalent* to having contributed  $g_i + \gamma$ .<sup>14</sup> This implies that (i) the optimal punishment points in Stage 3, based on the equalization of monetary payoffs after punishment between the  $n'$  punishers and the punished, needs to be adjusted to  $p_{ij} = \frac{(g_i + \gamma) - g_j}{n' - c}$  and (ii) the condition that ensures that enforcers find it optimal to punish is the same as in FS (condition (b) above). This follows because the decision to punish or not only depends on preferences and the cost of punishment, as contributions and acquisition costs are sunk in Stage 3.

Armed with these two concepts, we perform an equilibrium analysis of the game. An important quantity will be  $\hat{n}$ , defined as the minimum number of conditionally cooperative enforcers required to sustain cooperation. Assuming that the inequity averse players in the group are symmetric, e.g. they share the same values of  $\alpha$  and  $\beta$ , then from condition (b) we find

$$\hat{n} = \text{int} \left[ 1 + \frac{(1 + \alpha)(n - 1) - \frac{\alpha}{c}}{\alpha + \beta} \right] + 1$$

Note that  $n'$  and  $\hat{n}$  need not be equal. For some parameter values, for example, if  $c < \frac{\alpha}{(1 + \alpha)(n - 1)}$ , the numerator of the fraction in the definition of  $\hat{n}$  is negative and therefore  $\hat{n} = 1$ , i.e., one enforcer is enough to make the punishment threat credible, even if there are  $n' > \hat{n}$  such enforcers in the group.

We first show that for every subgame perfect Nash equilibrium outcome in FS with symmetric contribution level  $g \in [0, y]$ , there exists an equivalent equilibrium in our setting. By equivalent we mean a perfect equilibrium with full participation and the same contribution level  $g$ . Consider the following strategy:

- In stage 1, participate, i.e. provide the sanctioning institution.
- In stage 2, if all players participated, then contribute  $g$ . Otherwise, contribute zero if you did participate; if you did not participate, contribute zero if the number of conditionally cooperative enforcers who are participants is smaller than  $\hat{n}$ ; otherwise contribute  $\gamma$ .
- In stage 3, if all players participated and contributed  $g$ , then do not punish. If all players participated but there were deviations in the contribution stage, then play according to FS.<sup>15</sup> If participation is less than full, and there were no deviations in the contribution stage, then

<sup>14</sup> But it is *not* equivalent in terms of the level of the public good provided by the group.

<sup>15</sup> See Proposition 5 in FS, and its proof.

do not punish if you participated; if there were deviations in the contribution stage, play according to FS.

The novel part of this strategy is in Stage 2. Deviations from the prescribed participation behaviour are deterred by threatening to play the “worst” subgame perfect Nash equilibrium of the FS setting: participants contribute nothing; non-participants contribute  $\gamma$  if the number of conditionally cooperative enforcers with the punishing mechanism is at least  $\hat{n}$  and contribute nothing otherwise. The reason the contribution decision needs to be conditional on the participation decision and the number of enforcers who have participated is that in the case that the punishment threat is credible, non-participants will contribute  $\gamma$  to avoid being punished.

This strategy is subgame perfect because it follows a FS subgame perfect equilibrium strategy to avoid deviations in stages 2 and 3 and a FS subgame perfect equilibrium strategy that implements the lowest contribution to enforce participation decisions in stage 1.

**Lemma 1.** *For every subgame perfect Nash equilibrium in the FS setting, there is an outcome equivalent subgame perfect Nash equilibrium with full participation in the endogenous costly provision game.*

In our setting, there also exist subgame perfect Nash equilibria with less than full provision of the sanctioning institution. To see this, consider an adaptation of the previous strategy that starts by having a subset  $n'$  of conditional cooperative enforcers participating in Stage 1 and the remaining players not participating. Then, the strategy follows by requiring a contribution of  $g + \gamma$  for those not participating and a contribution  $g$  for participants. Deviations in the participation stage would be punished by resorting to the worst Nash Equilibrium of the contribution stage (as in the previous strategy) while deviations in the contribution stage would be punished by sending the corresponding punishment points (as in the previous strategy).

**Lemma 2.** *The endogenous costly provision game also has subgame perfect Nash equilibria with less than full participation,  $\hat{n} \leq n' < n$ . In this case, the profile of contributions is asymmetric if  $\gamma > 0$ : participants contribute  $g \in [0, y - \gamma]$  while non participants contribute  $g + \gamma$ .*

One might argue that a focal strategy is the natural separation between inequity averse and selfish players: inequity averse players acquire the punishment technology while selfish players do not acquire it but increase their contribution over the level of participating players to avoid

any pecuniary advantage that would be punished. But this focal equilibrium does not maximise the group payoffs if the number of conditionally cooperative enforcers in the group exceeds the threshold  $\hat{n}$ . The most efficient equilibrium would be the one in which the number of participants is kept to the minimum  $\hat{n}$ , maximising the level of contributions to the public good at minimal costs.

**Lemma 3.** *The subgame perfect Nash equilibrium that maximises the group payoff is the asymmetric equilibria where  $\hat{n}$  conditionally cooperative enforcers participate and contribute  $y-\gamma$  and non-participants contribute  $y$ .*

### 3.2 Hypotheses

The above analysis provides us with some guidance to analyse the data from our experiment. Notice, in the first place, that the model based on FS has multiple equilibria, both in terms of the number of participants providing the sanctioning institution and of the contribution levels. Precise predictions thus require the application of some refinements. We start with predictions on participation, i.e., provision of the sanctioning institution. For the parameter values in the experiment,  $n = 4$  and  $c = 1/3$ , the numerator in the definition of  $\hat{n}$  (the minimum number of enforcers required to successfully implement cooperation) is positive,  $(1 + \alpha)3 - 3\alpha = 3$ , implying that one enforcer is not enough in the experiment to sustain a cooperative outcome. As the case where the cost of providing the sanctioning mechanism is zero can be considered to be analogous to the standard punishment mechanism, and given that players successfully coordinate in StdPun, we hypothesize successful coordination in CTP0, which requires at least two enforcers.

**Hypothesis 1.** *The number of players providing the sanctioning institution in CTP0 is at least two.*

When acquisition is costly, there are two arguments that point in the direction of lower participation in CTP1 than in CTP0. The first one is efficiency which, according to Lemma 3, requires the minimization of the number of players providing the sanctioning institution in CTP1. The second is that a positive acquisition cost changes the out of equilibrium payoffs – because it affects the punishment points sent in Stage 3 – and this has been shown to adversely affect the equilibrium selection process (see discussion below).



**Hypothesis 2.** *The number of players providing the sanctioning institution in CTP1 is smaller than in CTP0.*

We do not expect instant coordination on a subgame perfect equilibrium, but an adjustment process towards one, if any. In this process, off-equilibrium play will play an important role. In this respect, the existence of a positive acquisition cost  $\gamma$  implies that the optimal number of punishment points to non-participants  $p_{ij} = \frac{(g_i + \gamma) - g_j}{n' - c}$  is increasing in  $\gamma$ .

**Hypothesis 3.** *For those providing the sanctioning institution, the mean level of punishment imposed on group members is higher in CTP1 than in StdPun or in CTP0.*

Turning to the targets of punishment, previous studies suggest that for punishment to increase group contributions, it must be targeted effectively at low contributors to “encourage” higher contribution levels (see, for instance, Fehr and Gächter, 2000). In line with these studies, we expect that those with lower than average contributions to the public good will be targeted for punishment.

Also, there are out-of-equilibrium profiles that trigger *rational* “antisocial” punishment in the sense of punishment sent by a player to another player who has contributed more. This will happen when  $\gamma > g_j - g_i > 0$ . This rational antisocial punishment will be more frequent, i.e., in more contribution profiles, the larger the acquisition cost  $\gamma$ . This suggests that, for a given level of punishment, some of that punishment will be diverted away from low contributors and directed towards above average contributors in CTP1.

**Hypothesis 4.** *Punishment is targeted at below average contributors in all punishment treatments. While punishment of above average contributors cannot be ruled out, punishment of low contributors is crowded out by anti-social punishment only in CTP1.*

An issue is the contribution levels to which play will converge. In this respect, a zero acquisition cost should not present any additional challenge to players over the standard FS setting (recall lemma 1). Hence, we expect successful coordination in CTP0, similar to StdPun.

**Hypothesis 5.** *Contribution levels are similar in StdPun and in CTP0. Further they are both significantly higher than in VCM.*

Lemma 2 poses some challenges to players on the coordination problem when the acquisition cost is positive. Now, there are two classes of subgame perfect Nash equilibria: symmetric

equilibria in which all players participate and contribute the same amount, and asymmetric equilibria. In the FS game, the criteria of symmetry and efficiency select full contribution as the unique prediction of the game. In our case, symmetry and efficiency do not go hand in hand, rendering the task of coordinating on an equilibrium more complex. The introduction of an acquisition cost increases the size of punishment *and* antisocial punishment, thus changing payoffs off the equilibrium path (in equilibrium there is no punishment). In coordination games, changes in off-equilibrium payoffs are known to affect the equilibrium selection process (see Cooper et al. 1990). Further, as shown in Rand et al. (2010), and as seen in Hermann et al. (2008), the presence of significant anti-social punishment can prevent the evolution of cooperation. There are thus two behavioural channels through which the existence of an acquisition cost negatively affects the selection of the equilibrium with the highest contribution level.

**Hypothesis 6.** *Contributions in CTP1 are lower than in StdPun and in CTP0.*

Nailing down the level to which contributions converge in our game is a complex issue as there are many different equilibria. The problem is exacerbated when the game is repeated, even over a finite number of periods as the complexity of the strategies grows exponentially. It is well known that if the stage game has a multiplicity of equilibria, outcomes that do not correspond to any equilibrium of the stage game can nevertheless be observed in the repeated version.

A final word about the minimum number of enforcers  $\hat{n}$  required to successfully implement cooperation. As shown, at least two enforcers are required to successfully implement cooperation in our treatments. This means that if there is only one enforcer, the punishment institution is not credible and the corresponding contribution level is 0. Therefore, the sanctioning institution will not support positive contribution levels.

**Hypothesis 7.** *One enforcer will not credibly enforce the sanctioning institution and therefore no cooperative outcome can be supported.*

In summary, we hypothesise that the exogenously provided sanctioning institution (StdPun) will be effective in raising contributions over levels observed in VCM. In addition, we hypothesise that the institution is provided to a lower extent when its provision is costly. We expect that effective sanctioning institutions will be provided by at least two individuals in CTP0 and that they will be as effective in CTP0 as in StdPun. However, we do not

expect the emergence of effective sanctioning institutions that can successfully raise cooperation in CTP1.<sup>16</sup>

#### 4. Results

The presentation of results is organised around the testing of Hypotheses 1 through 7. However, we also present additional results that are related to the repeated nature of the decision setting. Unless otherwise stated, Mann-Whitney (hereafter MW) tests are used to make comparisons across treatments.<sup>17</sup> Because subjects did not have information about other groups, each four-person group represents an independent decision-making unit. For these tests, an observation is thus the mean (averaged over all 20 rounds) per-round variable (e.g. contribution, punishment or earnings) by each group in a treatment. The results are organized around the initial treatments. Section 4.5 reports the results from the CTP5 treatment that was conducted as a test of robustness of the initial treatments.

##### 4.1 Institution Formation by Individuals

Figure 1 presents the mean number of members choosing to provide the sanctioning institution (providers) over time in the two CTP treatments. Aggregating across rounds, Table 2 presents the distribution of groups according to the number of participants in the punishment stage.

Figure 1 and Table 2 show that the average number of members providing the sanctioning institution in CTP0 is larger than 2. While there are very few rounds (2.31%) with zero participants, full provision is observed in only 10% of rounds. The mean number of participants in CTP0 is 2.41, i.e., about 60% of group members consistently provide the sanctioning institution. We therefore find support for Hypothesis 1.

**Result 1:** *The opportunity to choose to provide the sanctioning institution leads to an average participation rate greater than 2 in CTP0.*

---

<sup>16</sup> We do not present hypotheses on differences in efficiency across treatments. These depend on the magnitudes of punishment used relative to the increases in contributions. Moreover, punishment has been shown to lead to a clear efficiency increase only in very long repeated decision settings (Gächter et al. 2008).

<sup>17</sup> The results are robust to t-tests. We rely on the nonparametric tests to avoid making assumptions about the distribution of our data. All t-tests are available from the authors on request.

As Figure 1 shows, after the initial decision rounds, the average number of participants in the punishment stage is consistently lower in CTP1 than in CTP0. Table 2 shows that there is a shift in the distribution towards the lower end in CTP1 relative to CTP0. As shown, there are very few rounds with four participants in CTP1. In particular, there are 4 participants in only 4 percent of all decision rounds and 3 participants in only 13 percent of all decision rounds. The mean number of participants per-round is 1.60 in CTP1 and this is significantly lower than 2.41 in CTP0 ( $p = 0.0040$ ,  $n = 13$ ).<sup>18</sup>

To examine persistence at the individual level within groups, individuals are ranked in each group by the number of rounds in which they chose to participate in the punishment stage (1 = individual with the largest number of participation rounds in the group, 4 = individual with the smallest number of participation rounds in the group). Figure 2 presents the average number of rounds in which individuals of each rank chose to provide the institution.

As Figure 2 shows, individuals in the first three ranks choose to participate in the punishment stage in a greater number of rounds in CTP0 than in CTP1. MW tests confirm that the differences are significant ( $p = 0.0177$ ,  $0.0002$  and  $0.0179$  for ranks 1, 2 and 3 respectively). Thus, we find support for Hypothesis 2.

**Result 2:** *Average participation is lower in CTP1 than in CTP0. Individuals provide the punishment institution more persistently in CTP0 than in CTP1.*

In summary, in almost all decision rounds, fewer than four individuals choose to provide the sanctioning institution in both CTP treatments. Further, the introduction of a positive acquisition cost, though negligible, significantly reduces provision.

## 4.2 Use of the Sanctioning Institution

### 4.2.1 Amount of punishment used

Figure 3 (a) presents, across decision rounds, the mean frequency with which those providing the sanctioning institution assign punishment to others in their groups and Figure 3 (b) presents the mean amount of punishment assigned by those providing the sanctioning institution. We adopt the convention that the number of players providing the institution in StdPun is four.

---

<sup>18</sup> The mean is also significantly lower than four ( $p = 0.0000$ ).

Conditional on providing the institution, Table 3 presents the mean frequency of punishment and mean “per-capita” punishment imposed by those in the punishment role in each of the punishment treatments. In addition, it also presents mean aggregate punishment imposed at the group level in each treatment.

Both Figure 3(a) and Table 3 show that, those providing the institution are more likely to use punishment in CTP1 than in StdPun or in CTP0. MW tests confirm that the frequency of punishment is significantly higher in CTP1 than in StdPun ( $p = 0.0350$ ) and in CTP0 ( $p = 0.0502$ ). However, the difference between StdPun and CTP0 ( $p = 0.4626$ ) is not significant.

Figure 3(b) and Table 3 show a similar pattern for the mean amount of punishment used by those providing the institution, i.e., per-capita punishment. MW tests show that per-capita punishment is significantly greater in CTP1 than in StdPun ( $p = 0.0377$ ), but that the difference between StdPun and CTP0 is not significant ( $p = 0.2767$ ). However, MW tests also show that the difference between CTP1 and CTP0 is not significant ( $p = 0.1278$ ). Thus, we find mixed support for Hypothesis 3.<sup>19</sup>

**Result 3:** *Frequency of punishment by individuals providing the institution is significantly higher in CTP1 than in StdPun and somewhat higher than in CTP0. "Per-capita" punishment levels are significantly higher in CTP1 than in StdPun, but similar in CTP0 and CTP1.*

#### 4.2.2 Targeting of Punishment

When analysing punishment behaviour, we follow Fehr and Gächter (2000) and Sefton et al. (2007) and focus on deviations of an individual’s contribution from the average contribution of the others in the group. Figure 4 (a) shows the observed frequency with which an individual receives punishment in a round when the deviation of their contribution from the average contribution of the other three members of their group in that round is negative and when it is non-negative. Conditional on being punished, Figure 4 (b) shows the mean punishment received by individuals in a round as a function of their deviation in that round.

---

<sup>19</sup> While we do not have a hypothesis on aggregate punishment at the group level, Table 3 shows that aggregate punishment is highest in StdPun and is lowest in CTP0. However, the combination of lower participation rates and higher per-capita punishment in CTP1 renders all paired comparisons between treatments statistically insignificant (MW  $p > 0.50$  in all cases).

Figure 4 (a) shows that across all punishment treatments, those with negative deviations are punished more frequently than are those with non-negative deviations. Based on Sign-rank tests, the difference in the frequency of being punished between negative and non-negative deviations is significant in StdPun ( $p = 0.0029$ ) and CTP0 ( $p = 0.0019$ ), but is not significant in CTP1 ( $p = 0.1239$ ). Figure 4 (b) yields a similar result. In all cases, those with negative deviations receive more punishment than do those with non-negative deviations. Based on Sign-rank tests, the difference in absolute punishment received between negative and non-negative deviations is significant in StdPun ( $p = 0.0218$ ) and CTP0 ( $p = 0.0033$ ), but not in CTP1 ( $p = 0.3465$ ).<sup>20, 21</sup>

Similar to previous studies examining sanctioning institutions, the results reported above indicate that negative deviations are targeted for punishment in all punishment treatments. However, the frequency and amount of anti-social punishment in CTP1 is similar to that of punishment directed towards those with below average contributions. We thus find support for Hypothesis 4.<sup>22</sup>

**Result 4:** *In StdPun and in CTP0, negative deviations are punished more severely and more often than are positive deviations. In CTP1, however, the difference in frequency and intensity of punishment between negative and positive deviations is not significant.*

Rand et al. (2010) show that significant anti-social punishment can lead to negative reactions and the prevalence of “spiteful defectors”. In order to more fully understand differences in CTP0 and CTP1 with regard to choosing to participate in the punishment stage, we estimate individual level Probit regressions where the dependent variable is 1 if the individual chose to participate in the punishment stage in the round and is zero otherwise. The independent variables are a dummy for participation in the previous round, the lagged (absolute) deviation of the individual’s contribution from the average contribution of the others in the group and round dummies. To investigate if “blind revenge” or “anti-social behaviour” is a factor (Ostrom

---

<sup>20</sup> For these tests, an observation is the difference between the average (over all 20 rounds) punishment, or frequency of punishment, received by those with negative deviations and those with non-negative deviations in each group in a treatment. The number of observations in each treatment is thus equal to the number of independent groups in that treatment (see Table 1). Sign-rank tests are used to examine if this difference is statistically different from zero.

<sup>21</sup> The result is robust to finer partitions of the range of negative and non-negative deviations and to regression analysis.

<sup>22</sup> While the relevant test for Hypothesis 4 requires a comparison of differences *within* each treatment as reported above, we also compared the frequencies and means of punishment received by those with negative and positive deviations across treatments. Those with negative deviations are punished more frequently in StdPun than in CTP0 ( $p = 0.0043$ ) and in CTP1 ( $p = 0.0393$ ). All other differences are not significant.

et al., 1992, and Hermann et al., 2008), we also include the amount of punishment received by the individual in the previous round and the number of *other* participants in the punishment stage in the previous round.<sup>23</sup> To further check if revenge plays a role if received punishment was “anti-social” or “pro-social”, we run separate regressions for non-negative and negative lagged deviations. The results of this analysis are presented in Table 4. For all regressions, we report robust standard errors clustered on independent groups. For the sake of brevity, the coefficients of the round dummies are not reported.

The regressions suggest that there is strong path dependence in both CTP treatments in regard to participation; subjects who participate in one round are more likely to participate in the next round. However, the amount of punishment received in a round is a strong predictor of participation in punishment in the following round *only* in CTP1.<sup>24</sup> Further, this is the case whether a player’s contribution was below or above the average contribution level of others in the previous round.<sup>25</sup>

**Result 4a:** *Those who are punished are more likely to choose to participate in the punishment stage in the next round in CTP1, but not in CTP0.*

Result 4a is complementary to Result 4, which showed that there is significant anti-social punishment only in CTP1. Result 4a suggests that those choosing to participate in the punishment stage in CTP1 may have a greater tendency toward blind revenge or spite, targeting high contributors. The combination of these two results suggests that, in CTP1, the punishment of low contributors is crowded out by punishment targeted at those with positive deviations, leading to less effective use of punishment in increasing group contributions.

There is evidence from the field that a fear of retaliation prevents people from using punishment in the first place (Balafoutas and Nikiforakis, 2012 and Balafoutas et al., 2014). While retaliation seems to be a significant driver of punishment in CTP1, there is no evidence that

---

<sup>23</sup> When senders of punishment can be identified, there is evidence that players engage in targeted counter-punishment even in the same period (Denant-Boemont et al., 2007, Nikiforakis, 2008 and Nikiforakis et al., 2012) or in targeted revenge in the next period (Leibbrandt et al., 2015). In our setting, subjects were unable to identify who they received punishment from. Hence, we restrict attention to “blind” revenge.

<sup>24</sup> This result also holds when the independent variable is a dummy for receiving *any* punishment rather than the amount of punishment received. These regressions are presented in Table B1 in Appendix B.

<sup>25</sup> A possible explanation is that subjects find punishment to be a “hostile” act in CTP1, regardless of whether it is ‘pro-social’ or ‘anti-social’ punishment. Subjects could take it more personally if someone else has *paid* for the right to punish them, and thus set out to take (blind) revenge. Unfortunately, our data do not allow a test of this conjecture.

punishment is lower in CTP1 or that it declines over time. Indeed, Result 3 suggests that punishment use is higher in CTP1 than in StdPun and/or in CTP0. Results 4 and 4a suggest that subjects in CTP1 engage in “blind feuds” that lead to cycles of counter-punishment (Nikiforakis, 2008 and Nikiforakis et al., 2012) *over multiple rounds* of the game, resulting in counter-productive use of the sanctioning institution.<sup>26</sup>

However, we do not find evidence consistent with such multi-period feuds when participation in the sanctioning institution is costless; revenge or retaliation are not found to be strongly correlated with sanctioning in CTP0. Why then do we observe evidence of feuds only in CTP1? One possible reason is that the simple introduction of the monetary cost to the right to punish changes the attitude of participants toward other group members’ contribution and sanctioning decisions.<sup>27</sup>

#### 4.3 Effectiveness of Endogenously Provided Sanctioning Institutions

Figures 5(a) and 5(b) show the evolution of mean group contributions and earnings (both measured in tokens) over time. Since the initial endowment in each round was 20 tokens per individual (80 for the group) and all costs were paid out of this in all treatments, differences in earnings across treatments directly capture differences in efficiencies across treatments. Table 5 presents summary statistics of per-round group contributions and earnings.

Focusing first on contributions to the group fund, in all treatments mean contributions start at approximately 40 tokens (50% of the group’s endowment). Thereafter, contributions in the VCM and StdPun treatments follow a pattern similar to other studies examining these treatments (see, for instance, Fehr and Gächter, 2000). In VCM, they steadily decline over the course of the game to about 15 tokens (about 20% of endowment). In StdPun, they rise to around 60 tokens (75% of endowment) by round 5 and stay at that level throughout the rest of the game. The trajectory of contributions in CTP0 is very similar to that in StdPun.

MW tests support the observations made above. Compared to VCM, group contributions are significantly higher in StdPun ( $p = 0.0056$ ) and in CTP0 ( $p = 0.0004$ ). However, contributions

---

<sup>26</sup> Related to feuds is the literature on vendetta games, Abbink and Hermann (2009) find that fear of retaliation reduces money burning in groups of four, thus reducing the emergence of ‘pointless vendettas’. Bolle et al. (2014), however, find widespread retaliatory vendettas in stealing games in groups of two, leading to significant efficiency losses.

<sup>27</sup> There is evidence that the introduction of money reduces social distance among players and increases opportunistic behaviour (Vohs et al., 2006).



in StdPun and in CTP0 are not significantly different from each other ( $p = 0.4146$ ). We thus find support for Hypothesis 5.

**Result 5:** *Averaging across all 20 rounds, aggregate contributions are similar in StdPun and in CTP0. Moreover, they are both higher than in VCM.*

Mean contributions in CTP1 start similar to those in the other punishment treatments. They begin to rise in the first 2-3 rounds. While contributions in StdPun and CTP0 continue to rise, in CTP1 they then remain relatively flat throughout the game, above those in VCM but below those in the other two punishment treatments. However, they are closer to levels observed in VCM than in the other two punishment treatments. MW tests show that group contributions in CTP1 are not significantly different from those in VCM ( $p = 0.1069$ ) and that they are significantly lower than in both StdPun ( $p = 0.0296$ ) and in CTP0 ( $p = 0.0171$ ). Thus we also find support for Hypothesis 6.

**Result 6:** *Averaging across all 20 rounds, group contributions in CTP1 are significantly lower than in StdPun and in CTP0. Moreover, they are not significantly different from contributions in VCM.*

We next investigate to what extent increases in contributions to the group fund are linked to those who persistently choose to provide the sanctioning institution in their groups. Figure 6 presents mean contributions (over all 20 rounds) of individuals in each participation rank, as defined above (see Figure 2).<sup>28</sup>

The figure suggests that there is no difference in individual contributions by participation rank in CTP0. This is confirmed by an OLS regression (reported in Table B2 in Appendix B) where the dependent variable is an individual's mean contribution over all 20 rounds and the independent variables are dummies for participation rank within the group (excluded category: rank 4). None of the rank dummies is significant at the 10% level. In CTP1, the figure suggests that average contributions do not differ across the last three ranks but the average contribution of individuals with rank 1 is higher than that of the rest. However, an OLS regression (in Table B2 in Appendix B) shows that this difference is not significant. As above, none of the rank

---

<sup>28</sup> We do not present time trends of contributions of providers and non-providers. This is because an individual can be a provider in some rounds and non-provider in others. Calculating aggregate contributions by providers would thus involve potentially a different set of players in each round. Hence, we calculate separate averages for each individual in a group.

dummies is significant.<sup>29</sup> It thus appears that, in both CTP treatments, group contributions do not differ between those who participate persistently and those who do not.

However, there is a difference between the two CTP treatments in contributions by participation rank. Figure 6 also shows that mean individual contributions are higher in CTP0 than in CTP1 for each participation rank. MW tests show that this difference is not significant for rank 1 individuals ( $p = 0.2087$ ) but is significant for each of the other three ranks ( $p = 0.0096$ ,  $0.0129$  and  $0.0647$  for ranks 2, 3 and 4 respectively).<sup>30</sup> Thus providers of the sanctioning institution are more effective at raising contribution levels across ranks in CTP0 than in CTP1.

**Result 6a:** *Within each CTP treatment, group contributions are not significantly different between those that choose more often to provide the sanctioning institution and those that do less so. However, the contributions of those who provide the sanctioning institution are higher in CTP0 than in CTP1.*

While we do not have a hypothesis on group earnings or efficiencies, we nevertheless look at earnings ex-post. When comparing earnings across treatments, we account for the costs of punishment in the three treatments that allow players to punish each other. Figure 5(b) implies that these costs are substantial in the initial few rounds of the game. In the first five rounds, earnings in VCM are the highest while there is no discernible difference across the punishment treatments. In the remainder of the decision rounds, group earnings are lowest in CTP1 and are highest in CTP0. There is no systematic difference between earnings in VCM and earnings in StdPun. Further, they both lie in between earnings in the two CTP treatments. This is evident from the mean earnings in Table 5 as well. Mann-Whitney tests show that, across all 20 rounds the only pairwise comparison with a significant difference is the one between CTP0 and CTP1. In particular, group earnings are significantly higher in CTP0 than in CTP1 ( $p = 0.0129$ ).<sup>31</sup>

---

<sup>29</sup> In the OLS regressions in both treatments, the constant is positive and significant and is equal to the mean contribution of the rank 4 individual presented in Figure 6. The result is robust to individual-level panel random effects regressions (not reported) that include the above independent variables and lagged contributions and round dummies.

<sup>30</sup> Similar to Figure 6, we also created two figures of individual distributions with each rank in CTP0 and CTP1. The first was distributions of mean absolute individual contributions, while the second was distributions of mean individual deviations from the average contribution of the other three members of the group. Neither figure suggests any differences between treatments in any of the four ranks.

<sup>31</sup> Focusing on the last 10 rounds, earnings in CTP0 are significantly higher than in CTP1 ( $p = 0.0019$ ) and VCM ( $p = 0.0053$ ), but not significantly different than in StdPun.

**Result 6b:** *Averaging across all 20 rounds, mean group earnings in the three punishment treatments are very similar to earnings in VCM. Earnings in CTP1, however, are significantly lower than in CTP0.*

Thus, we find that the sanctioning institutions provided by individuals in CTP0 are as effective as when there is universal and exogenous participation in the sanctioning institution, i.e., in StdPun. However, the sanctioning institutions that emerge endogenously in CTP1 are not effective at raising contributions to the public good. The use of the sanctioning institution that emerges in CTP0 outperforms that in CTP1 in terms of both contributions *and* efficiency.

#### **4.4 Contributions: Level and Persistence of the Sanctioning Institution**

The previous results show that a smaller number of members provide the sanctioning institution in CTP1 compared to CTP0 and that contributions are lower in CTP1 than in CTP0. In terms of group outcomes, the question becomes to what extent contribution levels vary with the number and persistence of participants in the punishment stage. To examine this issue, Figure 7 presents mean contributions of groups according to the average number (over 20 rounds) of participants in the sanctioning institution. Recall, in StdPun, the number of participants is four in every round since all players automatically enter the punishment stage and is zero in the VCM treatment. The horizontal lines for these two cases represent reference points for average contributions.

Figure 7 provides evidence that group contributions increase with the average number of players persistently providing the sanctioning institution. Group level panel regressions (reported in Table B3 in Appendix B) of group contributions on lagged contributions, the lagged amount of punishment used and the number of providers confirm the positive relationship in both CTP treatments.

Importantly, Figure 7 also provides evidence that group contributions in CTP0 are as high as in StdPun when at least 2 participants provide the sanctioning institution. MW tests confirm that mean contributions in groups with at least two providers in CTP0 ( $n = 10$ ) are not significantly different from group contributions in StdPun ( $p = 0.5097$ ). However, mean contributions in groups with fewer than two providers in CTP1 ( $n = 9$ ) are significantly lower than in StdPun ( $p = 0.0330$ ).

Note that comparisons between StdPun and groups with fewer than two providers in CTP0 or groups with at least two providers in CTP1 are not very meaningful due to the small number of observations in the CTP treatments. In CTP0, 77% of the groups have at least two persistent providers of the sanctioning institution while in CTP1, only 31% have at least two persistent providers. A proportions test shows that this difference between the two CTP treatments is significant ( $n = 13$  groups in each; two-sided  $p = 0.0183$ ).

**Result 7:** *Groups contributions in CTP0 and CTP1 are significantly and positively correlated with the number of group members providing the sanctioning institution. Moreover, group contributions in CTP0 are as high as in StdPun when at least two players consistently provide the institution.*

#### 4.5 The Effect of a Non-negligible Participation Fee

In CTP1, the participation fee of one token is negligible. Nevertheless, it is still positive. Thus, CTP1 allows us to test if the mere presence of a positive price impacts behaviour and efficiency. The results discussed above show that when subjects have to pay to acquire the punishment technology, the sanctioning institution is provided to a significantly lower extent. Moreover, the use of the institution and group outcomes are very different. The experiments discussed in this section examine the extent to which the impact of requiring a positive price to provide the sanctioning institution varies with the magnitude of the price. In an additional treatment (CTP5), the participation fee was raised to 5 tokens, one-quarter of a subject's per-round endowment. All other details were identical to those in CTP1.

Three sessions of CTP5 were conducted at EssexLab, each lasting approximately 55 minutes. The average earnings of a subject in this treatment was £12 including a £2.50 show-up fee. Data were collected on 11 independent groups.

Figure 8 presents the mean number of members providing the sanctioning institution across rounds in CTP1 and CTP5. As an additional reference, the figure also presents the information for CTP0. As shown, the monotonicity argument continues to hold; participation in the punishment stage steadily declines as the cost of participation rises. The mean number of providers in a round in CTP5 was 0.36 (st. dev. = 0.24). Based on an MW test, this is significantly lower than the average of 1.6 participants in CTP1 ( $p = 0.0005$ ).

Based on MW tests, the following additional conclusions can be drawn. As a result of the low rates of provision of the sanctioning institution, the average punishment used by a group in a round was also significantly lower in CTP5 than in CTP1 (2.42 vs. 5.35,  $p = 0.0238$ ). However, mean individual punishment by providers in CTP5 was significantly higher than in CTP1 (5.55 vs. 3.26,  $p = 0.0051$ ) and in CTP0 (5.55 vs. 2.16,  $p = 0.0001$ ). Finally, as in CTP1, the punishment of individuals with negative deviations from the group's average contribution are crowded out by punishment of those with positive deviations in CTP5. While individuals with negative deviations from the group average were punished more often than were those with non-negative deviations (in 17% of the instances vs 15%), this difference is not significant ( $p = 0.9292$ ). Conditional on receiving any punishment, those with negative deviations received 3.9 tokens in punishment while those with non-negative deviations received 3.66 tokens in punishment. Once again, this difference is not significant ( $p = 0.7221$ ).<sup>32</sup>

Figures 9 (a) and (b) show, respectively, the mean group contributions and mean group earnings over time in both treatments. As shown, the patterns of contributions and earnings are quite similar across decision rounds. The mean per-round group contribution and group earnings in CTP5 were 40.10 tokens and 108.16 tokens respectively (the corresponding standard deviations were 19.01 and 23.73). Neither of these is significantly different from those observed in CTP1.<sup>33</sup>

The results from CTP5 lend support to the overall robustness of the effect of a positive price for providing the sanctioning institution. In particular, patterns in contribution, participation and aggregate punishment decisions closely mirror those observed in CTP1. Further, they also support Hypothesis 3 that an increase in the acquisition cost increases “per-capita” punishment by providers and Hypothesis 4 that a positive provision cost leads to significant anti-social punishment.

---

<sup>32</sup> See Figure 4 for the corresponding values for CTP1. The unit of observation is the difference in average (over all rounds) punishment, or frequency of punishment, received by those with negative deviations and those with non-negative deviations in each group.  $n = 11$  for both tests.

<sup>33</sup> For the corresponding values in CTP1, see Table 5. As before, MW tests are used to compare contributions and earnings between treatments. The unit of observation is the mean (over all 20 rounds) per-round group contribution or earning for each group in a treatment. Thus,  $n = 13$  in CTP1 and  $n = 11$  in CTP5.

**Result 8:** *The effects of a positive acquisition cost on the formation of sanctioning institutions by individuals and their effectiveness in raising contributions are robust to non-negligible acquisition costs.*

## 5. Conclusions

The decentralised sanctioning institution is one of the most widely studied solutions to the free-rider problem in public goods games (for a recent review of the literature see Chaudhuri, 2011). Given the second-order free-riding problem (Yamagishi, 1986), an important issue is the emergence of the institution. Unlike previous studies that have explored exogenous provision of the institution or *group* choice as to whether to adopt the institution, this study explores the willingness of individuals to unilaterally provide and make use of the sanctioning institution.

We find that individuals are willing to unilaterally provide the institution in their groups. However, the level of utilization and effectiveness of the institution varies importantly as to whether the provision cost is zero. When provision is costless group members consistently provide the institution for themselves, although not at an individual rate of 100%. Further, in this case, the sanctioning institution is as effective as when it is exogenously and universally provided. Punishment is effectively targeted at low contributors, raising contributions to the public good.

We also find, however, that if provision of the institution requires the payment of even a minimal fixed cost, provision and effectiveness of the sanctioning institution decline. In the presence of a negligible monetary cost, the number of individuals who are willing to provide the sanctioning institution is insufficient to raise cooperation. Further, in this case, revenge appears to be a greater reason for individuals choosing to participate in the sanctioning institution. This motive renders punishment ineffective as punishment of low contributors is crowded out to a greater degree by punishment of high contributors resulting from blind revenge.

How do these experimental results, including ours, inform us about behaviour in the field characterised by social dilemmas and punishment possibilities? It is our view that the stylized experimental settings for linear public goods games with and without punishment are but idealised constructs that help us understand basic incentives and behaviour. The most critical idealisations concern the free flow of noiseless information on contributions among all

participants (for the VCM) and the free availability of the punishment technology (for the StdPun). Although one can find real life situations that meet these conditions, for example norm violation in a public space and punishment by telling violators off (Nikiforakis, 2008), in many others these idealisations are compromised and it is here where our experimental study is informative.

The most direct real world examples to come to one's mind with no free availability of the punishment technology is that of *vigilantes* or *security patrols* - individuals who make an investment to organize themselves and who, in the absence of any sanctioning institution, voluntarily and unilaterally assume the task of punishing norm violators. This example implicitly includes one key feature of our design: the public announcement of norm enforcement. This would correspond to vigilantes publicly patrolling the streets with the instruments of punishment in their hands (e.g. guns or radio transmission devices to alert the police). There are other cases, less dramatic, also characterised by involved individuals stepping out and making public commitments to enforce a norm. An example is the case of honour codes at schools or colleges with students publicly announcing that they would turn in classmates for cheating/plagiarism.<sup>34</sup> Similarly, in the political arena, there are cases of political parties publicly announcing that they would prosecute any other party who engaged in corrupt activities. In Spain, this was the case for UPYD – Union, People and Democracy - for the last ten years.<sup>35</sup> Our results are informative of behaviour in such situations.

In a seminal work, Ostrom et al. (1992) established that “self-governance is possible” in groups faced with a social dilemma. Since then, other work has examined the effectiveness of exogenously provided sanctioning institutions across a diverse set of treatment conditions and in situations where the institution is adopted at the group level through voting mechanisms. Our study adds to this literature.

We find that individuals acting *unilaterally* may be able to provide a form of “governance” in their groups and raise cooperation levels. However, we also find that the sanctioning institution provided by individuals can be fragile. The results reported here point to the important role that

---

<sup>34</sup> It is the personal experience of one of the authors at a graduate school in the UK in the mid nineties, where some written exams were conducted without any official monitoring. It was common for students to warn others not to cheat in the exam or they would be reported to the instructor immediately.

<sup>35</sup> Unfortunately, this party has run out of money (and electoral support) and has issued a public statement some months ago that they would stop prosecuting other parties.

participation costs may play in the willingness of individuals to provide a sanctioning institution and, importantly, in how it is used.

### **Acknowledgements**

The authors thank Maria Bigoni, Jeff Carpenter, David Cooper, Simon Gächter, Ron Harstad, Martin Kocher, Andreas Leibbrandt, Miguel Ángel Meléndez-Jiménez, Daniele Nosenzo, Charles Noussair, Ragan Petrie, Martin Sefton, Sigrid Suetens, Matthias Sutter, Jean-Robert Tyran, Erte Xiao, the Associate Editor, two anonymous referees and seminar participants at the 2014 ESA North American meetings, the 2013 Southern Economic Association Meetings and the University of Málaga for their helpful comments, suggestions and advice. Funding from the Spanish Ministry of Economy and Competitiveness (project ECO2014-52345-P), the School of Economics at the University of East Anglia and EssexLab is gratefully acknowledged.



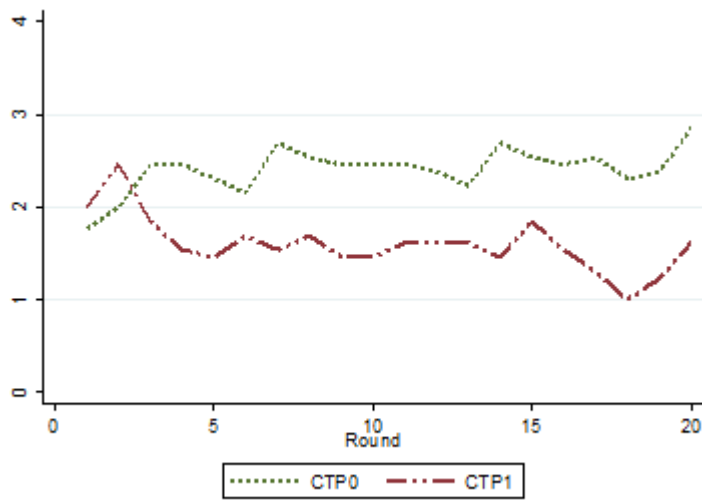
## References

- Abbink, Klaus and Benedikt Herrmann (2009) “Pointless Vendettas”, *CBESS Working Paper 09-10*.
- Anderson, Christopher M. and Louis Putterman (2006) “Do non-strategic sanctions obey the law of demand? The demand for punishment in the voluntary contributions mechanism”, *Games and Economic Behavior*, 54(1), 1-24.
- Balafoutas, Loukas and Nikos Nikiforakis (2012) “Norm enforcement in the city: A natural field experiment”, *European Economic Review*, 56(8), 1773-1785.
- Balafoutas, Loukas, Nikos Nikiforakis and Bettina Rockenbach (2014) “Direct and indirect punishment among strangers in the field”, *Proceedings of the National Academy of Sciences*, 111(45), 15924-15927.
- Bolle, Friedel, Jonathan H.W. Tan and Daniel John Zizzo (2014) “Vendettas”, *American Economic Journal: Microeconomics*, 6(2), 93-130.
- Carpenter, Jeffrey (2007) “The demand for punishment”, *Journal of Economic Behavior and Organization*, 62(4), 522-542.
- Carpenter, Jeffrey, Shachar Kariv and Andrew Schotter (2012) “Network Architecture, Cooperation and Punishment in Public Good Games”, *Review of Economic Design*, 95(1), 1-26.
- Chaudhuri, Ananish (2011) “Sustaining cooperation in laboratory public goods experiments: a selective survey of the literature”, *Experimental Economics*, 14(1), 47-83.
- Cinyabuguma, Matthias, Talbot Page and Louis Putterman (2005) “Cooperation under the threat of expulsion in a public goods experiment”, *Journal of Public Economics*, 89(8), 1421-1435.
- Cooper, Russell W., Douglas V. DeJong, Robert Forsythe and Thomas W. Ross (1990) “Selection Criteria in Coordination Games: Some Experimental Results”, *American Economic Review*, 80(1), 218-233.
- Croson, Rachel, Enrique Fatas, Tibor Neugebauer and Antonio J. Morales (2014) “Excludability: A laboratory study on forced ranking in team production”, *Working Paper*.
- Dannenbergh, Astrid, Andreas Lange and Bodo Strum (2010) “On the Formation of Coalitions to Provide Public Goods – Experimental Evidence from the Lab”, *ZEW Discussion Paper No. 10-037*.
- Denant-Boemont, Laurent, David Masclet and Charles N. Noussair (2007) “Punishment, counterpunishment and sanction enforcement in a social dilemma experiment”. *Economic Theory*, 33(1), 145-167.

- Ertan, Arhan, Talbot Page and Louis Putterman (2009) “Who to punish? Individual decisions and majority rule in mitigating the free rider problem”, *European Economic Review*, 53(5), 495-511.
- Fatas, Enrique, Miguel A. Meléndez-Jiménez and Hector Solaz (2010) “An experimental analysis of team production in networks”, *Experimental Economics*, 13(4), 399-411.
- Fehr, Ernst and Simon Gächter (2000) “Cooperation and Punishment in Public Goods Experiments”, *American Economic Review*, 90(4), 980-994.
- Fehr, Ernst and Klaus M. Schmidt (1999) “A Theory of Fairness, Competition and Cooperation”, *Quarterly Journal of Economics*, 114(3), 817-868.
- Fischbacher, Urs (2007) “z-Tree: Zurich toolbox for ready-made economic experiments”, *Experimental Economics*, 10(2), 171-178.
- Gächter, Simon, Elke Renner and Martin Sefton (2008) “The Long-Run Benefits of Punishment”, *Science*, 322(5907), 1510.
- Gerber, Anke, Jakob Neitzel and Philipp C. Wichardt (2013) “Minimum participation rules for the provision of public goods”, *European Economic Review*, 64, 209-222.
- Gürek, Özgür, Bernd Irlenbusch and Bettina Rockenbach (2006) “The Competitive Advantage of Sanctioning Institutions”, *Science*, 312(5770), 108-111.
- Herrmann, Benedikt, Christian Thöni and Simon Gächter (2008) “Antisocial punishment across societies”, *Science*, 319(5868), 1362-1367.
- Kamei, Kenju, Louis Putterman and Jean-Robert Tyran (2015) “State or nature? Endogenous formal versus informal institutions in the voluntary provision of public goods”, *Experimental Economics*, 18(1), 38-65.
- Kosfeld, Michael, Akira Okada and Arno Riedl (2009) “Institution Formation in Public Goods Games”, *American Economic Review*, 99(4), 1335-1355.
- Kube, Sebastian, Sebastian Schaube, Hannah Schidelberg-Hörisch and Elina Khachatryan (2015) “Institution Formation and Cooperation with Heterogeneous Agents”, *European Economic Review*, 78, 248-268.
- Leibbrandt, Andreas, Abhijit Ramalingam, Lauri Sääksvuori and James M. Walker (2015) “Incomplete punishment networks in public goods games: experimental evidence”, *Experimental Economics*, 18(1), 15-37.
- McEvoy, David M., Todd L. Cherry and John K. Stranlund (2011) “The Endogenous Formation of Coalitions to Provide Public Goods: Theory and Experimental Evidence”, *University of Massachusetts Amherst, Department of Resource Economics, Working Paper No. 2011-2*.

- Markussen, Thomas, Louis Putterman and Jean-Robert Tyran (2014) “Self-Organization for Collective Action: An Experimental Study of Voting on Sanction Regimes”, *Review of Economic Studies*, 81(1), 301-324.
- Masclet, David, Charles N. Noussair and Marie-Claire Villeval (2013) “Threat and Punishment in Public Good Experiments”, *Economic Inquiry*, 51(2), 1421-1441.
- Nikiforakis, Nikos (2008) “Punishment and counter-punishment in public good games: Can we really govern ourselves?”, *Journal of Public Economics*, 92(1-2), 91-112.
- Nikiforakis, Nikos, Charles N. Noussair and Tom Wilkening (2012) “Normative conflicts and feuds: The limits of self-government”, *Journal of Public Economics*, 96(9-10), 797-807.
- Oliver, Pamela (1980) “Rewards and Punishments as Selective Incentives for Collective Action: Theoretical Investigations”, *American Journal of Sociology*, 85(6), 1356-1375.
- Ostrom, Elinor, James Walker and Roy Gardner (1992) “Covenants with and without a Sword: Self-Governance is Possible”, *American Political Science Review*, 86(2), 404-417.
- Rand, David G., Joseph J. Armao IV, Mayuko Nakamaru and Hisashi Ohtsuki (2010) “Anti-social punishment can prevent the co-evolution of punishment and cooperation”, *Journal of Theoretical Biology*, 265(4), 624-632.
- Sefton, Martin, Robert Shupp and James M. Walker (2007) “The Effect of Rewards and Sanctions in Provision of Public Goods”, *Economic Inquiry*, 45(4), 671-690.
- Shampanier, Kristina, Nina Mazar and Dan Ariely (2007) “Zero as a Special Prize: The True Value of Free Products”, *Marketing Science*, 26(6), 742-757.
- Sutter, Matthias, Stefan Haigner, and Martin G. Kocher (2010) “Choosing the Carrot or the Stick? Endogenous Institutional Choice in Social Dilemma Situations”, *Review of Economic Studies*, 77(4), 1540–1566.
- Tralusen, Anne, Torsten Röhl and Manfred Milinski (2012) “An economic experiment reveals that humans prefer pool punishment to maintain the commons”, *Proceedings of the Royal Society B*, 279(1743), 3716-3721.
- van der Heijden, Eline, Jan Potters and Martin Sefton (2009) “Hierarchy and opportunism in teams”, *Journal of Economic Behavior and Organization*, 69(1), 39-50.
- Vohs, Kathleen D., Nicole L. Mead and Miranda R. Goode (2006) “The Psychological Consequences of Money.” *Science*, 314(5802), 1154-1156.
- Yamagishi, Toshio (1986) “The provision of a sanctioning system as a public good”, *Journal of Personality and Social Psychology*, 51(1), 110-116.
- Zhang, Boyu, Cong Li, Hannelore De Silva, Peter Bednarik and Karl Sigmund (2014) “The evolution of sanctioning institutions: an experimental approach to the social contract”, *Experimental Economics*, 17(2), 285-303.

**Figure 1. Mean number of providers of the sanctioning institution**



**Figure 2. Persistence in participation by individuals**

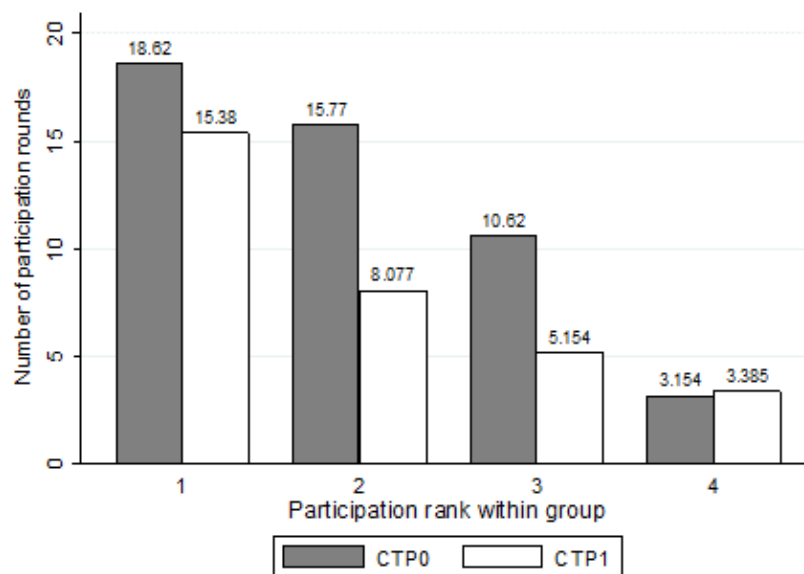


Figure 3. Mean frequency of punishment and mean punishment by providers

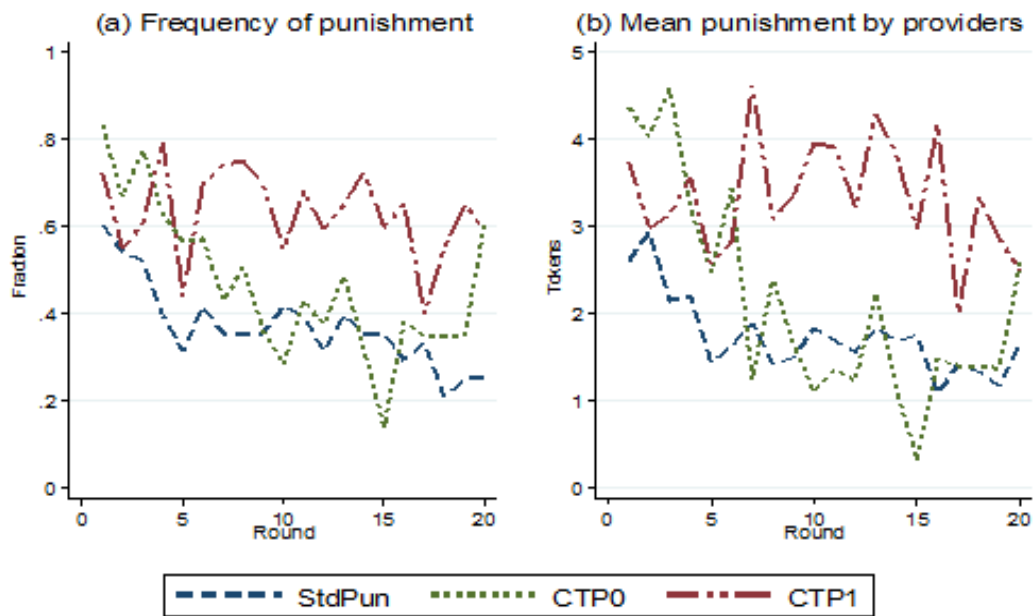


Figure 4. Frequency and amount of punishment received by individuals

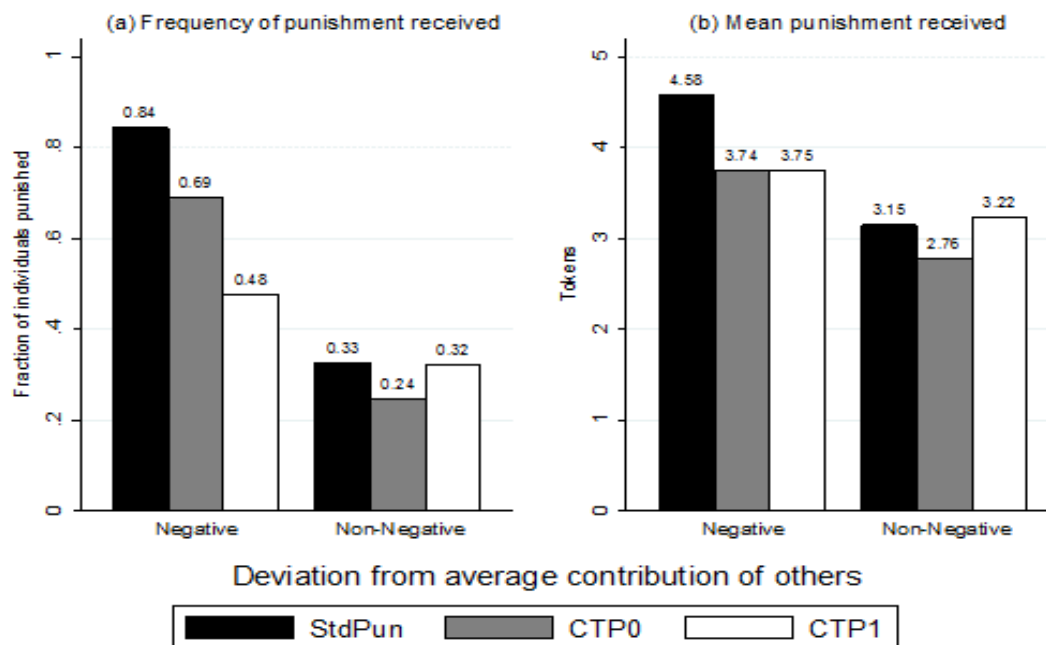
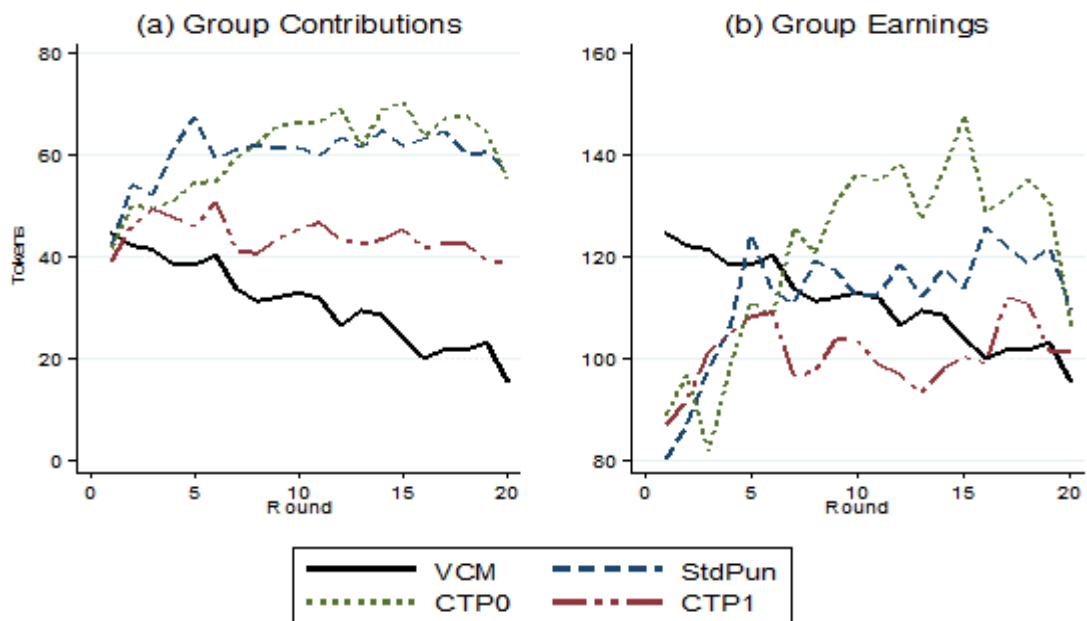
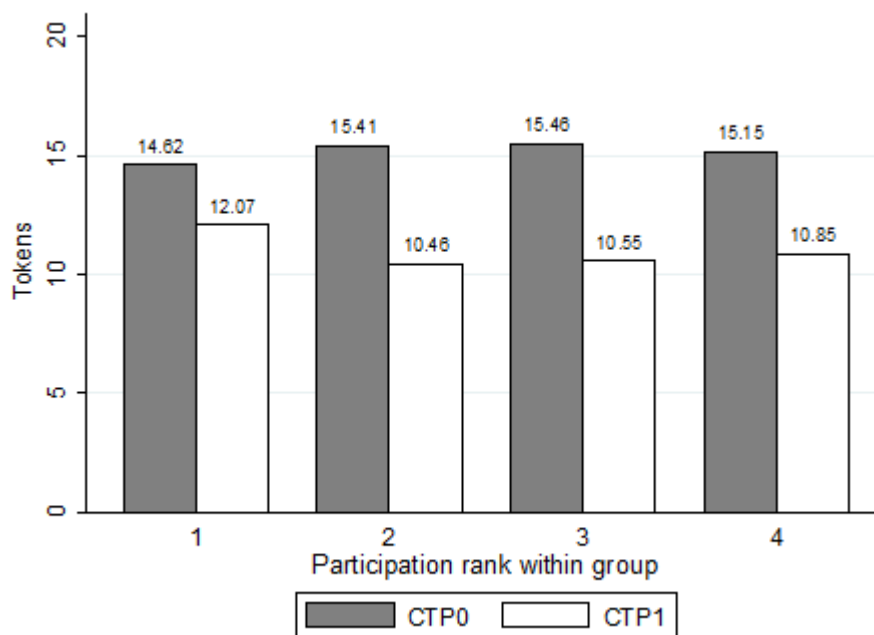


Figure 5. Mean Group Contributions and Earnings

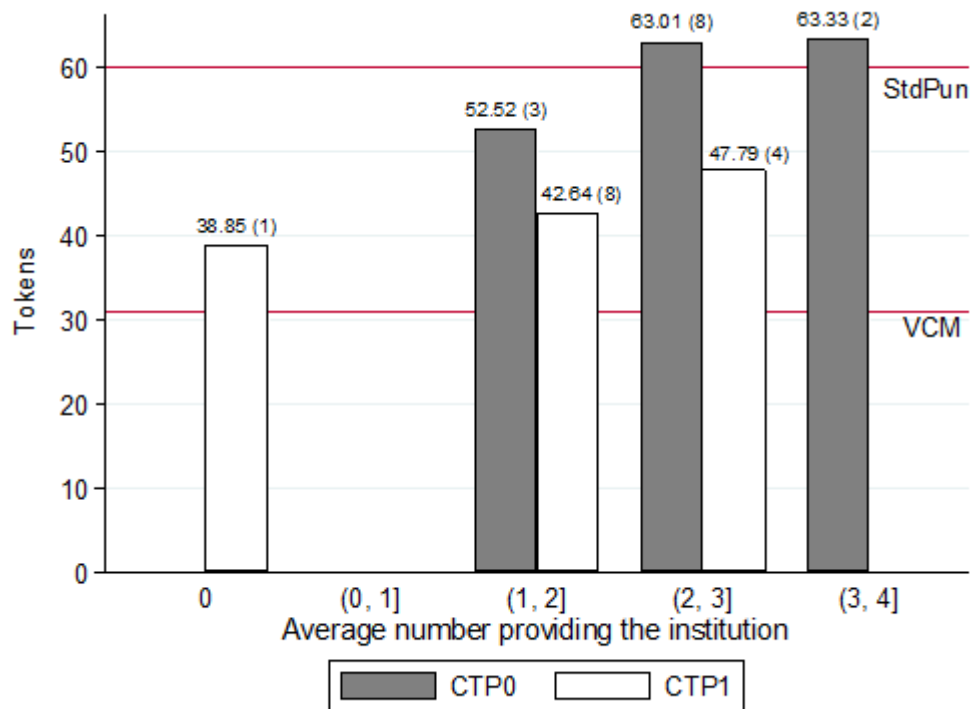


Note: Group earnings at the Nash equilibrium are 80 tokens.

Figure 6. Mean individual contributions over all 20 rounds by rank of participation in the sanctioning institution

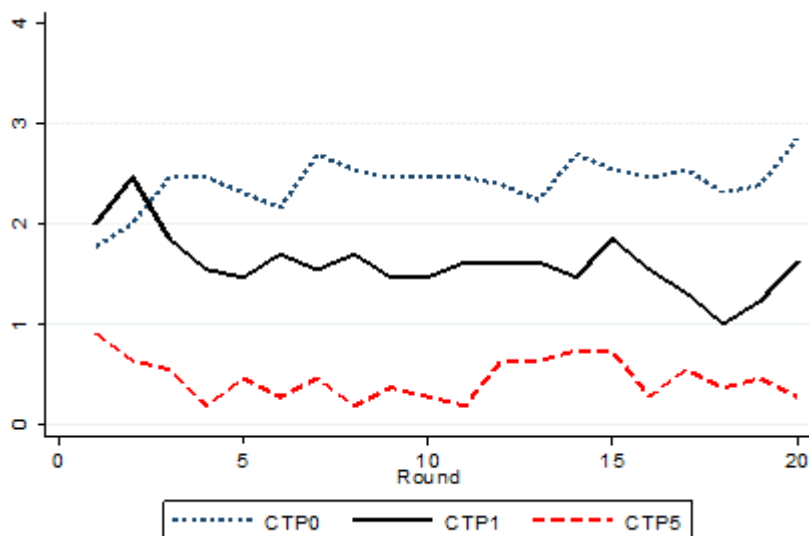


**Figure 7. Mean group contributions by number of members providing the institution**



**Note:** The number of participants is not always a whole number since it is an average over 20 rounds. Figures in parentheses are the number of groups in each category. There are 13 groups in each of the CTP treatments.

**Figure 8. Mean number of participants in the punishment stage – CTP treatments**



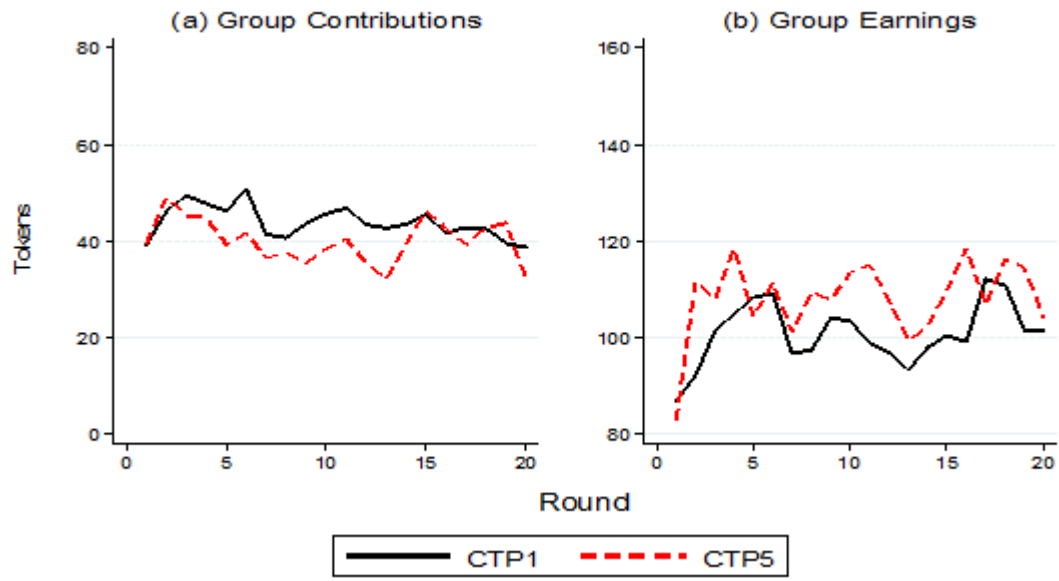
**Figure 9. Mean Group Contributions and Earnings in CTP1 and CTP5**



Table 1. Summary of initial treatments

Treatment	Punishment Opportunity	Participation in Punishment Stage	Punishment Participation Cost	Number of subjects (groups)
VCM	No	-	-	40 (10)
StdPun	Yes	All, automatically	-	48 (12)
CTP0	Yes	Only those who choose to in Stage 1	0 tokens	52 (13)
CTP1	Yes	Only those who choose to in Stage 1	1 token	52 (13)

Table 2. Distribution of the number of providers per group across rounds and mean participation in the punishment stage (all 20 rounds)<sup>36</sup>

# providers	% of rounds	
	CTP0	CTP1
0	2.31	16.15
1	15.38	28.85
2	31.54	38.08
3	40.77	12.69
4	10	4.23
<b>Mean Level</b>	2.41	1.60

<sup>36</sup> Figures in the table are percentages of groups in each category. Each group yields 20 observations, one for each round. Thus, each group could be in multiple categories. For instance, a group might have had 3 participants in punishment in round 10 but 2 participants in punishment in round 15.

**Table 3. Punishment Means (standard deviations) at the Group Level**

	Obs	Frequency of punishment	Per-capita Punishment	Aggregate Group Punishment
<b>StdPun</b>	12	0.371 (0.268)	1.739 (1.592)	6.958 (6.369)
<b>CTP0</b>	13	0.471 (0.208)	2.156 (1.403)	4.931 (2.791)
<b>CTP1</b>	12	0.624 (0.187)	3.258 (1.912)	5.35 (3.792)

**NOTE:** There are only 12 observations in CTP1 since there was one group where no one ever provided the sanctioning institution.

**Table 4. Determinants of participation in the punishment stage**

	Non-negative deviations <b>CTP 0</b>	lagged <b>CTP 1</b>	Negative deviations <b>CTP 0</b>	lagged <b>CTP 1</b>
Whether participated in the last round	2.240 <sup>***</sup> (0.183)	1.665 <sup>***</sup> (0.157)	1.892 <sup>***</sup> (0.218)	1.441 <sup>***</sup> (0.148)
Amount of punishment received in the last round	0.011 (0.041)	0.076 <sup>***</sup> (0.028)	-0.021 (0.051)	0.087 <sup>**</sup> (0.035)
Lagged <i>absolute</i> deviation from the average contribution of others	0.001 (0.020)	0.012 (0.023)	-0.012 (0.021)	-0.004 (0.018)
Number of <i>other</i> participants in the last round	-0.024 (0.105)	-0.055 (0.115)	0.003 (0.171)	0.029 (0.161)
Constant	-0.710 <sup>*</sup> (0.424)	-0.849 <sup>**</sup> (0.333)	-0.551 (0.496)	-1.625 <sup>***</sup> (0.343)
Observations	702	580	286	408

Dep. variable: = 1 if the choice was to participate in punishment stage and = 0 otherwise in each round. Std. errors clustered on independent groups in parentheses. Includes round dummies (not reported). \*\*\* - sig. at 1% level, \*\* - sig. at 5% level, \* - sig. at 10% level.

The higher number of observations in the regressions for non-negative lagged deviations indicates that there were more instances where individuals' contributions were at least as high as the average contribution of the others in the group. It does *not* indicate that there were more instances of 'anti-social' punishment than 'pro-social' punishment – the data includes observations where individuals received zero punishment.

**Table 5. Means (standard deviations) at the group level measured in tokens**

	<b>Obs</b>	<b>Contributions</b>	<b>Earnings</b>
<b>VCM</b>	10	31.01 (14.212)	111.01 (14.212)
<b>StdPun</b>	12	60.021 (20.472)	112.188 (41.311)
<b>CTP0</b>	13	60.639 (11.935)	120.915 (18.555)
<b>CTP1</b>	13	43.931 (18.103)	100.931 (19.897)