



International Institute for
Applied Systems Analysis
Schlossplatz 1
A-2361 Laxenburg, Austria

Tel: +43 2236 807 342
Fax: +43 2236 71313
E-mail: publications@iiasa.ac.at
Web: www.iiasa.ac.at

Interim Report

IR-13-078

**The evolution of sanctioning institutions:
An experimental approach to the social contract**

Boyu Zhang
Cong Li
Hannelore De Silva
Peter Bednarik
Karl Sigmund (ksigmund@iiasa.ac.at)

Approved by

Ulf Dieckmann
Director, Evolution and Ecology Program

June 2015

Interim Reports on work of the International Institute for Applied Systems Analysis receive only limited review. Views or opinions expressed herein do not necessarily represent those of the Institute, its National Member Organizations, or other organizations supporting the work.

The evolution of sanctioning institutions: an experimental approach to the social contract

Boyu Zhang ^{a,b}, Cong Li ^c, Hannelore De Silva ^d, Peter Bednarik ^e and Karl Sigmund ^{b,f,*}

^a School of Mathematical Sciences, Beijing Normal University, 100875 Beijing, China.

^b Faculty of Mathematics, University of Vienna, 1090 Vienna, Austria.

^c Key Laboratory of Animal Ecology and Conservation Biology, Centre for Computational Biology and Evolution, Institute of Zoology, Chinese Academy of Sciences, Beijing 100101, P. R. China.

^d Department of Finance, Accounting and Statistics, Vienna University for Economics and Business, 1190 Wien, Austria.

^e Courant Research Center Evolution of Social Behavior, University of Göttingen, D-37073 Göttingen, Germany.

^f International Institute for Applied Systems Analysis (IIASA), A-2361 Laxenburg, Austria.

*Author for correspondence: Telephone: +43 (0)1427750612

Fax: +43(0)142779506

E-mail: karl.sigmund@univie.ac.at

Abstract

A vast amount of empirical and theoretical research on public good games indicates that the threat of punishment can curb free-riding in human groups engaged in joint enterprises. Since punishment is often costly, however, this raises an issue of second-order free-riding: indeed, the sanctioning system itself is a common good which can be exploited. Most investigations, so far, considered peer punishment: players could impose fines on those who exploited them, at a cost to themselves. Only a minority considered so-called pool punishment. In this scenario, players contribute to a punishment pool before engaging in the joint enterprise, and without knowing who the free-riders will be. Theoretical investigations (Sigmund et al. 2010) have shown that peer punishment is more efficient, but pool punishment more stable. Social learning, i.e., the preferential imitation of successful strategies, should lead to pool punishment if sanctions are also imposed on second-order free-riders, but to peer punishment if they are not. Here we describe an economic experiment (the Mutual Aid game) which tests this prediction. We find that pool punishment only emerges if second-order free riders are punished, but that peer punishment is more stable than expected. Basically, our experiment shows that social learning can lead to a spontaneously emerging social contract, based on a sanctioning institution to overcome the free rider problem.

Keywords

Public goods, experiments, collective action, punishment, institution, social learning.

JEL codes

C73, C91, D03

1. Introduction

The role of punishment in boosting cooperation is a well-studied topic in experimental economics. However, most investigations deal with so-called peer-punishment (see, e.g., Fehr and Gächter 2000, 2002; Fehr and Rockenbach 2003; Fowler 2005; Henrich et al. 2006; Casari 2007; Sigmund 2007; Carpenter 2007; Egas and Riedl 2008; Gächter et al. 2008; Chaudhuri 2011). Typically, the players in a public good game are allowed to impose fines on others, at a cost to themselves. The threat of punishment can lead to considerable increases in the level of cooperation in the collective action. Many players are willing (and frequently even eager) to shoulder the costs of imposing fines on cheaters. Some subjects, however, abuse sanctioning opportunities by engaging in antisocial punishment which harms cooperators (Cinyabuguma et al. 2006; Denant-Boemont et al. 2007; Dreber et al. 2008; Hermann et al. 2008; Nikiforakis 2008).

In most aspects of everyday life, the task of punishing exploiters has eventually been taken over by institutions (Ostrom 2005; Guala 2012). In developed societies, peer-punishment is not only less common (Balafoutas and Nikiforakis 2012), but often explicitly forbidden. Under conditions of anarchy (e.g., Hobbes's 'state of nature'), individuals have to take punishment into their own hands, but in all better-regulated communities, individuals engage in some form of social contract by delegating punishment to institutions. How can we envisage this important step in social development?

Evidently, this question can be approached from many different angles. We use an economic experiment to test how individuals who want to coerce their group to cooperate decide between inflicting pro-social punishment directly or using the intermediary of an institution. We do not offer them opportunities for antisocial punishment, or the bribery of institutions, but rather concentrate on the idealized versions. The foremost problem, in such an experiment, is how to implement the sanctioning institution (Tyler and Degoey 1995; Casari and Luini 2009; Kosfeld et al. 2009; Andreoni and Gee 2012). Which is the essential feature distinguishing institutional from peer-punishment? Some argue that it is the delegation of punishment. However, individuals who want to exert personal revenge can recur to 'hiring a gun', and this would still count as peer-punishment (Van Vugt et al. 2009). A more pronounced difference is that sanctioning institutions are established in advance, and thus entail running costs even in the case that no one commits a punishable offense. A county would have to pay its sheriff even if nobody commits a crime. We tried to model this as 'pool-punishment' (for experimental papers, see Yamagishi 1986; Guillen et al. 2006; Kamei

et al. 2011; Markussen et al. 2011; Traulsen et al. 2012; for theory, see Sigmund et al. 2010, 2011). Players who want to use such a sanctioning tool have to pay a fee, even before the joint enterprise takes place, or at least before they are informed of its outcome, and thus before they know whether there will be any exploiters to punish. Pool punishers can be viewed as paying a tax towards a police. We note that instead of pool- or peer-punishment, some authors use the terms ‘formal’ and ‘informal’ sanctions (Kamei et al. 2011; Markussen et al. 2011).

In our experiment, we investigated 18 small groups, or ‘toy-communities’, of 12 to 14 players. Each such group played 50 rounds of a Mutual Aid game, isolated from the other groups. The Mutual Aid Game is a variant of a Public Goods game, where players do not obtain any return from their own contribution, and hence are faced with an even more pronounced social dilemma. Within each group, players could decide, before each round, whether to join a Mutual Aid game without punishment (NoPun), with peer-punishment (Peer), with pool-punishment (Pool) or not to participate in the game at all (No). These games were played separately, i.e., the outcome of one game did not affect the outcomes of the other games in the group. Players were anonymous, and prevented from communicating. All that players learned, after each round, was how many opted, in their group, for each alternative, and which payoff they obtained. They then could choose whether to opt for (NoPun), (Peer), (Pool) or (No) in the next round. We thus observed, in each toy community, whether social learning led to institutional punishment or not.

It is clear that if pro-social punishment works as desired, i.e., if it leads to all-out cooperation, then peer- punishment is more efficient than pool-punishment, since it entails no running costs. However, theoretical considerations (Sigmund et al. 2010, see relevant theory in section 2) imply that pool-punishment is more stable, provided that it is also directed at those participants in the game who do not contribute to the punishment pool. Indeed, if cooperation is achieved, i.e., if no one needs to be punished, then a peer-punisher cannot be distinguished from a non-punisher. This means that second-order free-riders (defined as those who contribute to the Mutual Aid, but not to the sanctions) cannot be spotted, and thus cannot be punished. By contrast, those who do not contribute to the punishment pool are just as visible as those who do not contribute to the Mutual Aid, and can be punished just as well. A system implementing this is highly immune against exploitation, but requires payment of some sort of a tax to maintain the punishment pool. We wanted to compare the attractiveness of first-order and second-order pool punishment against the background of the same peer punishment and no punishment alternatives.

In our experiment, a clear majority chose peer punishment in the first round. Most players switched to pool punishment in later rounds, but (as predicted by theory) only if punishment was also imposed on second-order free-riders. The experiment involved 238 first-year students from universities in Vienna. Interactions were anonymous. Players were randomly allocated to 18 groups of 12 to 14 players each, for the duration of 50 rounds. We implemented 2 treatments with 9 groups each: in the ‘second-order treatment’, players were offered a pool punishment game which sanctioned second-order free riders, and in the ‘first-order treatment’ a pool punishment game which did not. The former treatment led to the emergence of pool punishment in six out of the nine groups, the latter in none. Peer punishment slowly declined over rounds in both treatments. Roughly speaking, it was not displaced by pool punishment, but eroded gradually. Contributions to the Mutual Aid game were vastly more frequent in the treatment with second-order pool punishment.

In a nutshell, players were allowed to ‘vote with their feet’ (the expression seems to be due to Tiebout 1956), and most of them decided in favor of a sanctioning institution, but only if this institution coerced participants to contribute not merely to the Mutual Aid, but also to its own upkeep. Under this additional commitment, the institution was adopted by the group, in a kind of ‘social contract’ which was achieved without explicit communication or deliberation, and was based on social learning from the own experience and that of others.

In section 2, we describe the theoretical background, and in section 3, the experiment. In section 4, we display the results, and in section 5, we offer a discussion and conclusions. The instructions for the players and the detailed results of every group are contained in the supplementary information of the Online Resource.

2. Theoretical background: a choice of games

In this section, we introduce the types of games used in the experiment and briefly sketch some of the relevant theory from Sigmund et al. (2010).

First of all, let us consider the Mutual Aid game (MA game) of type (NoPun) (no punishment). There are m players in the group. They can decide whether or not to contribute an amount $c > 0$, knowing that this will be multiplied by $r > 1$ and divided among all *other* players in the group. If m_C is the number of those players who contribute, and m_D the number of those who don’t (with $m_C + m_D = m$), then the payoff for a contributor is

$$P^C = rc \frac{m_C - 1}{m - 1} - c$$

and that for a defector

$$P^D = rc \frac{m_C}{m - 1}.$$

Clearly, we always have $P^D > P^C$ (the difference is independent of m_C) and thus the dominant strategy is to refuse to contribute. If all players contribute, their payoff is $(r-1)c$, which is independent of group size m . In our experiment, $c=1$ monetary unit (MU), $r=3$ and $m \geq 2$ is variable. This game was first described in Wilson (1975), see also Sugden (1986), Yamagishi (1986), Fletcher and Zwick (2004), Sigmund (2010). It is very similar to the usual public good game (PG game), see e.g. Fehr and Gächter (2000). Whereas in the latter, an amount r/m of a player's contribution returns to the player, in the MA game players do not obtain any return from their own contribution. In Sigmund et al. 2010, 2011, the Mutual Aid game has been called the 'others only' variant of the Public Goods game. In our setup, the number of participants m can fluctuate. If m is smaller than r , the PG game is no longer a social dilemma: to contribute is the dominant strategy. Most PG games do explore the case when $m < r$, but we wanted to make it more difficult for cooperation to emerge: for the MA game, defection is always the dominant strategy. Note that if everyone contributes, the payoff is the same, namely $(r-1)c$, in both the MA and the PG game.

Now let us consider the MA game of type (Peer) (peer punishment). After deciding whether to contribute or not to the Mutual Aid in an MA game, players can punish defectors. We assume that only players who contribute can punish. If they do, they have to punish all defectors in the group. As the focus of our study is the decision between different types of pro-social punishment, this assumption makes the issue as clear-cut as possible.

Let us suppose that m_{Pe} is the number of players who contribute and punish those who do not contribute, m_C is the number of players who contribute, but do not punish, and m_D is the number of those who neither contribute nor punish, i.e., the defectors or free-riders. Thus $m_{Pe} + m_C + m_D = m$. Every peer punisher has to punish every free rider. Let β be the size of the fine that each non-contributor has to pay *per punisher*, and γ the fee each punisher has to pay for each non-contributor he or she punishes. The cost of punishment is uncertain as it depends on the number m_D of free-riders, but the fine-to-fee ratio is fixed to $\beta: \gamma$. Then we obtain as payoff values

$$P^C = rc \frac{m_C + m_{Pe} - 1}{m - 1} - c$$

$$P^{Pe} = rc \frac{m_C + m_{Pe} - 1}{m - 1} - c - \gamma m_D$$

$$P^D = rc \frac{m_C + m_{Pe}}{m - 1} - \beta m_{Pe}$$

There is no dominant strategy. The group optimum is obtained whenever $m_D = 0$. In this case, every player obtains $(r-1)c$, just as in the MA game of type (NoPun). Clearly, we have $P^C \geq P^{Pe}$ (with equality if and only if $m_D = 0$). The state when no one contributes is a strict Nash equilibrium. Other (non-strict) equilibria exist for $m_D = 0$ and $m_{Pe} \geq (c+\beta)/\beta$. In our experiment, $\beta = 1$ MU and $\gamma = 0.5$ MU so that states with two or more peer punishers, but no defector are also Nash equilibria.

Finally, let us consider the MA game of type (Pool) (pool punishment). Just as in the (Peer) game, we assume that only contributors can punish. At the same time that they decide whether to contribute to the Mutual Aid or not, they also decide whether to contribute to the punishment pool. There are m_C players who contribute to the Mutual Aid, but not to the punishment pool, m_{Po} players who contribute to both pools, and m_D players who contribute to neither pool (with $m_{Po} + m_C + m_D = m$). Pool punishers have to contribute an amount c to the Mutual Aid and an amount F to the punishment pool. The fine-to-fee ratio depends on the number of free-riders.

The (Pool) game is played in two variants. In the first-order variant, everyone who does not contribute to the Mutual Aid is fined by an amount Bm_{Po} , whereas in the second-order variant, everyone who does not contribute to the punishment pool (i.e., who does not contribute to both Mutual Aid and the sanctioning) is fined by that amount. The payoff values are

$$P^{Po} = rc \frac{m_C + m_{Po} - 1}{m - 1} - c - F$$

and in the first-order variant (Pool1)

$$P^C = rc \frac{m_C + m_{Po} - 1}{m - 1} - c$$

$$P^D = rc \frac{m_C + m_{Po}}{m - 1} - Bm_{Po}$$

resp. in the second-order variant (Pool2)

$$P^C = rc \frac{m_C + m_{P_o} - 1}{m - 1} - c - Bm_{P_o}$$

$$P^D = rc \frac{m_C + m_{P_o}}{m - 1} - Bm_{P_o}$$

In our experiment, we used $B=1$ MU and $F=0.5$ MU. In the first-order variant, we have $P^C > P^{P_o}$ so that $m_D = m$ is the only Nash equilibrium. In the second-order variant, $m_{P_o} = m$ is another equilibrium (as long as $c + F \leq B(m-1)$), which for our parameter values means that there are at least three punishers). This equilibrium is not efficient (i.e., pareto-optimal), since $m_C = m$ provides a higher per capita payoff and actually is the group optimum.

In each of these games, the all-defector state is a strict Nash equilibrium. In the context of evolutionary games, it represents a social trap. How, then, could cooperation emerge? In Sigmund et al. 2010, it was shown that if the game is optional, i.e., if players have also the possibility of abstaining from it, then cooperation based on peer- or pool punishment can emerge. Indeed, the participation in a Mutual Aid game (or a public good game) can be viewed as a risky enterprise, which only succeeds if enough co-players cooperate. Non-participants can be viewed as risk-averse players who, rather than engage in such an uncertain interaction, prefer to engage in some other activity whose payoff σ does not depend on what the others are doing. We assume that this payoff is somewhere between the payoff obtained if no one contributes to the MA game, and that obtained if everyone contributes to the common pool (and, in the Pool game, to the punishment pool). This means $0 < \sigma < (r-1)c$ and, in the Pool game, $0 < \sigma < (r-1)c - F$. Whenever defectors thrive and cooperation breaks down, the option of abstaining from the game yields a higher payoff and non-participants take over. This in turn gives players willing to invest in the Mutual Aid game the chance to re-establish cooperation. The possibility of abstaining from the game therefore offers an escape from the social trap. In the theoretical model of Sigmund et al. 2010, free riders took over if participation was not optional, but compulsory.

For the competition between peer-punishment and pool-punishment, Sigmund et al (2010, 2011) show, using arguments from evolutionary game theory, that in the first-order variant, peer-punishers prevail most of the time, but sometimes second-order free-riders invade. In this case, defectors and then non-participants take over before peer-punishment is reestablished. In contrast, in the second-order version, pool-punishers eventually establish a very stable regime, although it is less efficient.

3. The experiment

The 18 groups of 12 to 14 players (our ‘toy-communities’) were the independent sample points of our experiment. Players in different groups were not allowed to communicate with each other and interacted only within their group. The players were not told that the number of rounds was fixed beforehand at 50, so as to prevent end-round effects. In each round, players were given 3 MU and asked to choose one of three variants of the Mutual Aid (MA) games: (NoPun) MA without punishment; (Peer) MA with peer punishment; (Pool) MA with pool punishment. The players could also decide (No) not to participate in any of these games. Such non-participants received an additional 0.5 MU. The idea, here, was that when not participating in a joint enterprise, an individual can engage in some useful activity which does not depend on the decisions of others.

Once players had chosen one of these games, they played one round of the game they had chosen with those group-members who had chosen the same game. Players who opted for one of the games (NoPun), (Peer) or (Pool), but found no co-players to join them, were treated as non-participants (No), and received an additional 0.5 MU, independently of what the others did. Once the round was over, the players learned how many (in their group) had played (NoPun), (Peer), (Pool) or (No), how many in each game had chosen which strategy, and which payoff they had obtained. They could use this information to decide for which game to opt in the next round. Players did not learn about who did what, so there was no possibility to establish a reputation. Players knew that they would be paid immediately after the experiment, at a rate of 10 cents (euro) per MU, without having to give away their identity (as players) to their co-players or to the experimenters. The guaranteed minimal payoff was 10 euro.

Players participating in a MA game of type (NoPun) could decide whether or not to contribute 1 MU to the common pool, knowing that their contribution would be multiplied by 3 and divided equally among all *other* players in their game, irrespective of whether these co-players had contributed or not. Thus contributors did not benefit from their own contribution. If all cooperate, everyone gains 2 MU.

Players choosing to participate in an MA game of type (Peer) would first play an MA game as described above, and then, in a second stage of the same round, be shown the number of non-contributors (i.e., defectors) in their game. Contributors could then decide whether or not to punish these free-riders. The fine-to-fee ratio is fixed to 2:1 in (Peer), as used in Carpenter (2007) and Nikiforakis and Normann (2008). Each punisher would have to pay a

fee of 0.5 MU per defector, and each defector would have to pay a fine of 1 MU per punisher. Again, if all cooperate, everyone gains 2 MU.

Players participating in a MA game of type (Pool) had to choose between three options: (i) not to contribute anything, (ii) to contribute to the Mutual Aid (i.e., to pay 1 MU so that 3 MU would be shared among all other members who had chosen (Pool)), or (iii) to contribute to *both* the Mutual Aid and the punishment pool. This last alternative requires the players to pay 1 MU to the Mutual Aid and an additional 0.5 MU into the punishment pool. Thus if all cooperate, everyone gains 1.5 MU. This MA game was played in two variants, denoted as ‘first-order variant’ resp. ‘second-order variant’. In the first-order variant, players knew that everyone who had not contributed to the Mutual Aid would be fined 1MU per punisher. In the second-order variant, players knew that everyone who had not contributed to *both* Mutual Aid and punishment pool would be fined 1MU per punisher. Hence, in the second-order variant of game (Pool), the institution punishes all of these free-riders (i.e., those who opted for (i) or (ii)) irrespective of whether they contributed or not to the Mutual Aid game), while in the first-order variant (Pool1), second-order free-riders (those who opted for (ii)) were not punished. The fine to fee ratio can greatly vary, in this game, depending on the number of defectors and pool punishers. In groups 1-9 (with altogether 120 subjects), the game of type (Pool) was offered in the first-order variant, and in groups 10-18 (with 118 subjects) in the second-order variant.

We note that this is a complex game, without obvious money-maximizing strategies for the individuals choosing (Peer) and (Pool), since payoff depends on how many decide for the different options. In order to provide the players with an appreciation of the issues involved, they were given, at the start of the session, 25 practice rounds (see Online Resource). They knew that these rounds would not count towards their score and that groups would be reshuffled before the experiment started. More precisely, players were first given, via computer screen, a brief introduction into game (NoPun), then played five rounds of the game (NoPun). The same then happened with games (Peer) and (Pool). Finally, they all played 10 rounds with the option, in each round, to choose between the three games (NoPun), (Peer) (Pool), or (No), which meant to abstain from participation (exactly as later in the actual experiment). Thus players could familiarize themselves with their options, in the practice rounds, but were precluded from sharing their experiences through communication. Immediately after the practice rounds, the ‘toy communities’ were re-assembled randomly.

After each round, players were shown the payoffs for all strategies used in their group, and had 15 seconds to decide which game (NoPun), (Peer), (Pool) or (No) to join next. The tightness of the schedule and the complexity of the task provided a strong motivation to be guided by the size of the payoffs, i.e., to engage in social learning. We also did use loaded language in the instructions, for instance by calling punishment ‘punishment’. Since our main aim was to compare different pro-social sanctioning mechanisms, we felt justified in acknowledging the underlying, common intention to uphold norms of collaboration. In the same spirit, asocial punishment or revenge were not offered as options to our players. Moreover, to reduce the complexity, we avoided the issue of increasing or decreasing group returns (i.e., we assumed that the Mutual Aid was proportional to the number of contributors). Furthermore, we considered only pure strategies. Players could either make a full contribution or none, and either punish all exploiters or none. This need not correspond to real-life situations, but rather aims at eliciting clear responses from the candidates facing choices between pro-social punishment mechanisms. Admittedly, these experimental strictures limit the generalizability of the results.

4. Results: social learning of social control

In the actual experiment, we observed strong changes in behavior in most of the 18 groups, especially during the initial phase. 12 of the groups eventually settled down, in the sense that the same game was chosen by the majority for each of the last 10 rounds. Six of these groups settled down for pool punishment. All six belonged to the second-order treatment. In three groups playing the second-order treatment, and three groups playing the first-order treatment, players settled for peer punishment. The null hypothesis that pool punishment is equally likely in both treatments can be rejected with a significance of $p < 0.05$ ($n_1=9$, $n_2=9$, two-sided binomial sample test). Based on the theoretical model, we had indeed expected pool punishment to emerge in the second-order treatment only.

The average frequency of pool punishment increased during the first rounds, in the second-order treatment, and overtook the frequency of peer punishment. In fact, the initial frequencies of (NoPun), (Peer), (Pool) and (No), in the first-order treatment, corresponded closely to the initial frequencies in the second-order treatment, but then the frequencies evolved very differently (see Figure 1). Frequencies of peer punishers decreased in both treatments, but only slowly. Frequencies of pool punishment decreased in the first-order

treatment, but increased in the second-order treatment. (See Online Resources for the regression equations)

More precisely, in the first round of the second-order treatment, 55 per cent of players choose the peer punishment game and 36 per cent the pool punishment game. The initial frequencies in the first-order version were 56 per cent and 31 per cent, respectively. However, in the first-order treatment, both frequencies declined to reach 48 per cent and 19 per cent, respectively, by round 50. By contrast, the evolution in the second-order treatment reversed frequencies, so that after 50 rounds, 63 per cent of players opted for the pool punishment game but only 33 per cent for the peer punishment game (Figure 1b). This reversal took place in the first 20 rounds. The regression equation is $y=0.326+0.0146x$ (where y represents the frequency of pool-punishment and x the round), with coefficient of determination $R^2=0.9167$ (which measures how well the regression line represents the data) and $P\text{-value}<0.0001$. Obviously, players approached both first- and second-order treatments with similar expectations, but then underwent a very different learning experience.

If we average over all 50 rounds, we find a significant preference for peer punishment in the first-order treatment, and a less significant preference for pool punishment in the second-order treatment (Figure 3a). The latter treatment leads to a very pronounced cooperative behavior. Indeed, the frequency of contributions was significantly higher in the second-order treatment than in the first-order treatment (88.2 per cent vs. 48.9 per cent, Mann-Whitney U-test, $n_1=9$, $n_2=9$, $p=0.0373$), and it hardly changed over the 50 rounds (Figure 2b). We can see (Figure 3c and Online Resource) that average payoff values differ by little, but that peer punishment clearly yields the highest payoff in the first-order treatment, whereas it shares front rank with pool punishment, in the second-order treatment.

In the first-order treatment, peer punishment was preferred by a wide margin: game (Peer) was chosen in 55.6 per cent of all decisions, game (Pool) in 20.2 per cent, game (NoPun) in 11.7 per cent and non-participation (No) in 12.6 per cent (Figure 3a). A majority (62 per cent) of decisions in the peer punishment game (Peer) were to contribute to the Mutual Aid, but not to punish, only 12 percent were for punishment. (In both treatments, the average number of peer punishers in the peer games is decreasing with the number of defectors. See Table S4 in Online Resource). The payoff of these first-order defectors was higher than that of the punishers (4.636 MU vs. 4.1 MU, Mann-Whitney U-test, $n_1=9$, $n_2=9$, $p=0.077$). It is obvious that within any round, this has to hold, if some players defect; we see here that it also holds on average. The non-contributors in the peer punishment game earned

marginally more than the non-participants (3.61 MU vs 3.5 MU, the difference is not significant). All in all, 48.9 per cent of all decisions were in favor of contributing to the Mutual Aid, rather than defecting (35.6 per cent) or abstaining from the game (15.5 per cent). But the time evolution over 50 rounds tells a more pessimistic story (Figure 2a). Three-fourth of players cooperated in the first round but half of them gave up in later rounds. The regression equation is $y=0.669-0.0065x$ (where y represents the frequency of cooperation and x the round), with coefficient of determination $R^2=0.8812$ and P-value <0.0001 . Moreover, in the first-order pool punishment games, neither contributions to the Mutual Aid nor to the sanctioning took off. In particular, only a tiny fraction of the decisions (54 out of 1149) favored investing into the punishment pool.

In the second-order treatment, the preferences for the games change drastically (the hypothesis that the preferences for the four games are the same can be rejected, using a Chi-square test, P-value <0.0001). The game (Pool), was chosen in 54.1 per cent of all decisions, and almost always (namely, in 3155 of 3174 cases) was combined with a decision to actually contribute to the punishment pool. The peer punishment game (Peer) was chosen in 41 per cent of the decisions. Interestingly, players who chose the peer punishment game rarely decided to actually punish (only 9 per cent did), and the average payoff for those who actually engaged in peer punishment, 3.78 MU, was significantly less than that of second-order free-riders (4.77 MU, Mann-Whitney U-test, $n_1=9$, $n_2=9$, $p=0.0106$). But this minority of peer-punishers sufficed to keep free-riding in the Mutual Aid game down to 16 per cent. Few decisions (4.5 per cent) were in favor of the alternative (NoPun), i.e., joining a Mutual Aid game without punishment, and only 0.4 per cent were in (No). The average payoff for the peer punishment game (Peer) was insignificantly larger than for the pool-punishment game (Pool) (4.49 MU vs. 4.46 MU), but those who actually peer-punished had a significantly lower payoff than those who actually contributed to the punishment pool (3.78 MU vs. 4.49 MU, Mann-Whitney U-test, $n_1=9$, $n_2=9$, $p=0.004$).

The average payoff for those choosing a given game is almost the same for both first order and second order treatments, with one exception: the payoff for choosing pool punishment has substantially increased in the second-order treatment, because almost all players contributed to the Mutual Aid in the second-order treatment, but less than a third did so in the first-order version (Figure 3c).

It remains to show that players, in this experiment, were essentially guided by social learning, defined here as preferential copying of the strategies with the highest payoff

observed among other players in the same group. There were 2850 decisions to switch to another strategy. Of these, 2012 decisions (70.6 percent) were in favor of a strategy having currently a strictly better payoff in the group. There were 7446 decisions not to switch, but to repeat the former move. In 76.1 percent of the cases, the payoff was optimal. Conversely, when the payoff was optimal, 85.7 percent (i.e., 5666 out of 6615) of the decisions were not to switch.

The clearest form of social learning occurs when a player who switches adopts the strategy with currently highest payoff. This occurs indeed for 1700 of the 2850 switches. The frequency of players who adopt a strategy which is less than optimal decays exponentially in d , the difference between the current optimum and the current payoff of the strategy that the switching player will adopt in the next round. (Nonlinear regression $\nu = 0.4058 \times 0.2888^d$, correlation coefficient $R=0.8439$ and $P\text{-value}<0.001$, see Figure S3 in Online Resource). Needless to say, in the next round both the optimal payoff and the payoff of the newly adopted strategy may be different.

If we define a player as social learner when at least 90 percent of that player's decisions could be explained either as (a) switching to a strategy with currently higher payoff, or as (b) sticking with a strategy of currently highest payoff, then 78.6 percent of the players were social learners. We stress that the players had only 15 seconds between rounds, which hardly offered them time for strategic calculations.

5. Discussion and Conclusion

Coercion plays an essential role in overcoming social dilemmas. The corresponding line of reasoning goes back at least as far as Hobbes' 'Leviathan' from 1651, and the practical implementation can be traced throughout history. The selfish motivations endangering collective actions have to be suppressed by positive and negative incentives (Olson 1965; Boyd and Richerson, 1992; Andreoni et al., 2003; Rockenbach and Milinski, 2006). In particular, the threat of punishment curbs the temptation to free-ride, i.e., to exploit the contributions of others without offering an adequate return.

Institutions can be viewed as tools for providing incentives (Ostrom, 2005). It has been shown that even in small-scale societies far removed from 'Leviathan'-like states, grass-root institutions can deal, often efficiently, with the tasks of monitoring joint efforts and sanctioning defectors (Ostrom, 1990; Boehm, 2000; Henrich, 2006; Baldassarri and

Grossman, 2011). We wanted to test how players could opt for such a rudimentary institution, modeled as pool-punishment.

Our experiment is close in spirit and design to an experiment by Gürer et al. 2006. In that experiment, players were given the choice between a Public Good game with and one without peer punishment. The majority started with a clear preference for the treatment without punishment, but switched after a few rounds to the peer-punishment treatment, apparently guided by payoff considerations. Essentially, we kept the three-staged structure (choice of treatment, decision to contribute, decision to punish), but added pool punishment and non-participation as additional choices. (In contrast to the paper by Gürer et al. 2006, we did not allow for rewarding; a related endogenous choice between peer punishing and rewarding has been investigated by Sutter et al. (2010), who found that reward was often chosen.)

The option of pool punishment adds an important element, as it essentially provides the opportunity for a tacit social contract establishing a sanctioning institution. To our knowledge, this is the first experiment demonstrating that such a social contract can emerge through social learning based on comparing the (frequency dependent) payoff values of diverse options. ‘Social contract’ means that players can submit to a sanctioning authority, captured here as a ‘punishment pool’. By contrast, peer punishment corresponds to self-justice, and belongs to an anarchic ‘state of nature’. Both philosophers and experimental game theorists have shown that peer punishment can lead to an escalation of conflicts, i.e., to a ‘war of all against all’. For example, John Locke wrote in §126 of his ‘Two Treatises of Government’ from 1689 that in the state of nature, ‘... resistance [by defaulters] many times makes the punishment dangerous, and frequently destructive, to those who attempt it’. In a similar vein, experiments such as those by Denant-Boemont, Masclet, Noussair (2007), Nikiforakis (2008) or Nikiforakis and Engelmann (2011) show that peer punishment can invite counter-punishment and lead to costly feuds. We have deliberately excluded the possibility of counter-punishment, which threatens self-justice; we also have excluded the possibility of a corrupt authority, which threatens sanctioning institutions. The different punishment regimes offered in our experiment were, in this sense, idealized, pro-social versions.

The lesson for institution design is clear: pool punishment requires the sanctioning of second-order free-riders. The important role of second-order free-riding is well-known (Oliver 1980), and our experiment confirms it. In the second-order treatment, pool punishment effectively prohibits this possibility, whereas in the first-order treatment, it does not.

Apparently, pool-punishers notice that they are exploited, in the first-order treatment, and react against this breach in equity (Bolton and Ockenfels, 2006). Voting for the second-order treatment implies a higher commitment.

We now discuss several aspects of the experimental design which may limit the generalizability of our results.

We did allow for players to abstain from the game. Clearly, there exist joint enterprises or common resources from which one cannot abstain: the global climate is the best example. Such compulsory interactions do not belong to the class considered here, since we have allowed players to opt for non-participation. Nevertheless, it could well be that the main ‘efficiency vs. stability’ result still holds for compulsory games. It was for two reasons that we decided to consider only voluntary interactions in our experiment: first, because the theoretical results guiding our predictions were derived for this class of games only, and second because, in the course of the experiment, we sometimes (but rarely) encounter a player who is the only individual choosing a given Mutual Aid game of type (NoPun), (Peer) or (Pool). In this case, it is practical to assign them option (No), namely ‘non-participation’. This, incidentally, hardly affects the statistics.

The fact that we chose a Mutual Aid game, rather than the Public Good game, should not overly restrict the range of applications. We did so because the number of players was variable, and that the Public Good game is stops being a social dilemma if the number of players is smaller than the multiplication factor r . In contrast, a Mutual Aid game always is a social dilemma: this makes the issue of altruism vs. free-riding more clear-cut. As a real-world example of a Mutual Aid game, we refer to the ‘sick clubs’ or ‘friendly societies’ run by working men in nineteenth-century England (see Sugden, 1986). If one of the members fell ill, the others could contribute to his aid. Such joint enterprises were informal fore-runners of state insurance schemes. If we assume that in a given period, everyone is equally likely to fall ill, and that the others could decide between contributing a fixed amount or not, we obtain exactly the structure of a Mutual Aid game.

In our experimental design, we did not allow for punishment of non-punishers in the peer punishment game. The reason is twofold. On the one hand, theoretical models predict that it has no effect on the outcome (Sigmund et al, 2010). On the other hand, economic experiments have confirmed this in similar situations (Cinyabuguma et al. 2006; Kiyonari and Barclay 2008; Traulsen et al. 2012). We do not expect second-order peer punishment to affect the outcome.

We have reduced all individual decisions to choices between two, three or four alternatives. It would be interesting to investigate scenarios where players have a larger range of strategies, for instance by allowing them to choose between ten levels of contribution to the Mutual Aid, or different degrees of punishment. Similarly, we have proposed only one, extremely rudimentary form of institution. It is easy to think of better designs, for instance by allowing part, at least, of the unused funds to return to the players who have contributed to the punishment pool. We refrained from doing this, because we did not want to make it too easy for institutional punishment to emerge. The fact that as many groups ended up with peer- as with pool-punishment suggests that we succeeded in this ‘calibration’. Moreover, our experiment is already complex enough, and we feared to make it cognitively too demanding by adding more choices. As it was, the practice rounds needed to familiarize the players with their options took almost one hour (as long as the subsequent experiment).

Our main objective was to compare two different versions of pool punishment (rather than pool with peer). We note that there exist at least three experiments (independently conceived and in part not yet published) comparing pool with peer punishment, or ‘informal’ with ‘formal’ sanctions (Kamei et al. 2011; Markussen et al. 2011; Traulsen et al. 2012). Kamei et al. 2011 and Markussen et al. 2011 adopt a continuous version of first-order pool punishment, and Traulsen et al. 2012 consider both first- and second-order pool punishment (same as our experiment). In Markussen et al. 2011, fixed groups of five players play for 24 rounds, and can vote, at specific instants, between two different regimes (corresponding, in our setup, to decisions between (NoPun) and (Peer), (Peer) and (Pool), or (NoPun) and (PoolC)). In Kamei et al. 2011, the choice is between (Peer) and (Pool) with various parameters for the sanctions. Informal sanctioning does remarkably well, both from the viewpoint of frequency and payoff. (The experiment by Ertan et al. 2009 and the theoretical model by Boyd et al. 2010 confirm that peer punishment works well when players have an opportunity for coordinating.) In contrast, formal sanctions (which did not include second-order punishment) fared poorly.

The experiment by Traulsen, Röhl and Milinski 2012 investigates three treatments. In (a), players are offered the possibility of peer punishment, first 10 rounds without, then 15 rounds with second-order punishment. In (b), players are offered the possibility of pool punishment, again first without, then with second-order punishment. In (c), players could use both types of punishment in each round, by investing into a punishment pool before the public good interaction, and exerting peer punishment after the public good interaction. Again, this was first played without, and then with second-order punishment. In (a), switching from first-

to second-order punishment had very little effect; in (b) and (c), the switch strongly boosted pool punishment. In our experiment, players were not allowed to use both types of punishment simultaneously. They had to decide between one or the other (or none), thereby going separate ways, and building communities with different rules.

In all these experiments (including ours) peer punishment did reasonably well. This may in part be due to the higher efficiency in the optimal situation when all contribute to the Mutual Aid, or the Public Good. It may also be due to the fact that possibilities for retaliation and feuding were excluded in the experimental designs.

Since we wanted to favor conditions for social learning, we provided the players with information on the frequencies and payoffs obtained by the various strategies in their group (See Online Resource). However, we refrained from giving them opportunities to build up individual profiles, for instance reputations, or significant differences in resources. Needless to say, this does not imply that reputations or differences in resource holding power are irrelevant for the evolution of institutions. Similarly, we did not consider other-regarding preferences (Fehr and Schmidt 1999) or contests between groups, although such struggles played doubtlessly an important role in human evolution (Choi and Bowles 2007).

What are the roots of sanctioning institutions? Our players were given the choice between one type of peer and one type of pool punishment. Needless to say, such an approach cannot tell how such opportunities for sanctioning emerge. Cooperation has frequently arisen through biological evolution (Maynard Smith and Szathmary 1995), often via subtle mechanisms suppressing competition (Frank 1995), and there exist many examples of animals punishing each other (Clutton-Brock and Parker 1995). In particular, parents repress competition between their offspring, in many species, and it may be that this eventually led, in human populations, to institutionalized sanctioning. Offspring would simply have to remain with their parents (a costly option which provides some safety) rather than leave and defend their interests single-handedly. This fits with Jean Jacques Rousseau's claim, in 'Du Contrat Social' from 1762, that '*the family is the first model of a political society: its head is the image of the father*' (Book 1, Chapter 2). However, such a scenario was not addressed in our experiment. We presumed that the interactions are symmetric, and that all participants have equal resources. This better fits with Boehm 2000, who proposed another origin for sanctioning institutions. Its first instances, accordingly, were coalitions directed against alpha-males in the group, and the first social contract aimed at suppressing bullying behavior, in order to guarantee an egalitarian distribution of big-game meat.

It seems that institutions, once they have arisen, apply themselves to curb the vengeful and aggressive instincts fuelling peer-punishment. It would be interesting to explore this, both by modeling and by experiment. In our experimental setup, we have not allowed pool-punishers to sanction peer-punishers, or punished players to retaliate (Cinyabuguma et al. 2006; Nikiforakis 2008). We also excluded communication and deliberation, although theoretical models, field observations and experiments alike have stressed the importance of communication in sanctioning exploiters (Walker et al. 2000; Bochet et al. 2006; Ertan et al. 2009). If individuals can look for allies, or deliberate with their peers, stable systems of incentives can arise (Casari and Luini 2009; Ertan et al. 2009; Boyd et al. 2010). We aimed for a minimalistic scenario based on social learning, and showed that it can lead to the emergence of a rudimentary type of institutionalized coercion helping to overcome individuals' selfish preferences.

Acknowledgements

Our software was based on z-Tree, for which we thank Urs Fischbacher. We are grateful to Dirk Semmann, Jean Pierre Tyran, Matthias Sutter and Simon Gächter for comments, and to Stephan Aigner for help during the experiment. We acknowledge support from the Austrian Science Funds and the European Science Foundation through TECT I-104 G11 and Grant #RFP-12-21 from Foundational Questions in Evolutionary Biology Fund. We also thank for the support of the Chinese Scholarship Council (CSC).

References

- Andreoni, J. and Gee, L. L. (2012). Gun for hire: Delegated enforcement and peer punishment in public goods provision. *Journal of Public Economics*, 96, 1036-1046.
- Andreoni, J., Harbaugh, W. and Vesterlund, L. (2003). The carrot or the stick: rewards, punishments, and cooperation. *American Economic Review*, 93, 893-902.
- Balafoutas, L. and Nikiforakis, N. (2012) Norm Enforcement in the City: A Natural Field Experiment forthcoming *European Economic Review*.
- Baldassarri, D. and Grossman, G. (2011). Centralized sanctioning and legitimate authority promote cooperation in humans. *Proceedings of the National Academy of Sciences*, 108, 11023–11027.
- Bochet, O., Page, T. and Putterman, L. (2006). Communication and punishment in voluntary contribution experiments. *Journal of Economic Behavior and Organization*, 60, 11-26.
- Boehm, C. (2000) Conflict and the Evolution of Social Control, *Journal of Consciousness Studies* 7, 79-101.
- Bolton, G. E. and Ockenfels, A. (2006). ERC: a theory of equity, reciprocity, and competition. *American Economic Review*, 90, 166-93.
- Boyd, R., Gintis, H. and Bowles, S. (2010). Coordinated punishment of defectors sustains cooperation and can proliferate when rare. *Science*, 328, 617-620.
- Boyd, R., and Richerson, P. J. (1992). Punishment allows the evolution of cooperation (or anything Else) in sizeable groups. *Ethology and Sociobiology*, 113, 171-195.

- Carpenter, J. P. (2007). The demand for punishment. *Journal of Economic Behavior and Organization*, 62, 522–542.
- Casari, M. (2007). On the design of peer punishment experiments. *Experimental Economics*, 8, 107-115.
- Casari, M. and Luini, L. (2009). Cooperation under alternative punishment institutions: an experiment. *Journal of Economic Behavior and Organization*, 71, 273-282.
- Chaudhuri, A. (2011). Sustaining cooperation in laboratory public goods experiments: a selective survey of the literature. *Experimental Economics*, 14, 47-83.
- Choi, J-K. and Bowles, S. (2007). The coevolution of parochial altruism and war. *Science*, 318, 636-640.
- Cinyabuguma, M., Page, T. and Putterman, L. (2006). Can second-order punishment deter perverse punishment. *Experimental Economics*, 9, 265-279.
- Clutton-Brock, T. H. and Parker, G. A. (1995). Punishment in animal societies. *Nature*, 373, 209-216.
- Denant-Boemont, L., Masclet, D., Noussair, C. (2007) Punishment, Counterpunishment and Sanction Enforcement in a Social Dilemma Experiment, *Economic Theory* 33, 145-167.
- Dreber, A., Rand, D. G., Fudenberg, D. and Nowak, M. A. (2008). Winner don't punish. *Nature*, 452, 348–351.
- Egas, M., and Riedl, A. (2008). The economics of altruistic punishment and the maintenance of cooperation. *Proceedings of the Royal Society B: Biological Sciences*, 275, 871–878.
- Ertan, A., Page, T. and Putterman, L. (2009). Who to punish? Individual decisions and majority rule in mitigating the free rider problem. *European Economic Review*, 53, 495-511.
- Fehr, E. and Schmidt, K. M. (1999). A theory of fairness, competition, and cooperation. *Quarterly Journal of Economics*, 114, 817-868.
- Fehr, E. and Gächter, S. (2000). Cooperation and punishment in public goods experiments. *American Economic Review*, 90, 980-994.
- Fehr, E. and Gächter, S. (2002). Altruistic punishment in humans. *Nature*, 415, 137–140.
- Fehr, E. and Rockenbach, B. (2003). Detrimental effects of sanctions on human altruism. *Nature*, 422, 137–140.
- Fischbacher, U. (2007). Z-tree: Zurich toolbox for ready-made economic experiments. *Experimental Economics*, 10, 171–78.
- Fletcher, J.A. and Zwick, M. (2004). Strong altruism can evolve in randomly formed groups. *Journal of Theoretical Biology*, 228, 303-313.
- Fowler, J. H. (2005). Altruistic punishment and the origin of cooperation. *Proceedings of the National Academy of Sciences*, 102, 7047–7049.
- Frank, S. A. (1995). Mutual policing and repression of competition in the evolution of cooperative groups. *Nature*, 377, 520-522.
- Gächter, S., Renner, E. and Sefton, M. (2008). The long-run benefits of punishment. *Science*, 322, 1510–1512.
- Guala, F. (2012). Reciprocity: Weak or strong? What punishment experiments do (and do not) demonstrate. *Behavioral and Brain Sciences*, 35, 1-59.
- Guillen, P., Schwieren, C. and Staffiero, G. (2006). Why feed the Leviathan? *Public Choice*, 130, 115-128.
- Gürerk, O., Irlenbusch, B. and Rockenbach, B. (2006). The competitive advantage of sanctioning institutions. *Science*, 312, 108-111.
- Henrich, J. (2006). Cooperation, punishment, and the evolution of human institutions. *Science*, 312: 60–61.

- Henrich, J., McElreath, R., Barr, A., Ensminger, J., Barrett, C., Bolyanat, A., Cardenas, J. C., Gurven, M., Gwako, E., Henrich, N., Lesorogol, C., Marlowe, F., Tracer, D. and Ziker, J. (2006). Costly punishment across human societies. *Science*, 312, 1767-1770.
- Herrmann, B., Thoni, C. and Gächter, S. (2008). Antisocial punishment across societies. *Science*, 319, 1362–1367.
- Kamei, K., Putterman, L., and Tyran, J-R. (2011). State or nature? Formal vs. informal sanctioning in the voluntary provision of public goods. Discussion Papers 11-05, University of Copenhagen. Department of Economics.
- Kiyonari, T. and Barclay, P. (2008). Free riding may be thwarted by second-order reward rather than by punishment. *Journal of Personality and Social Psychology*, 95, 826-842.
- Kosfeld, M., Okada, A. and Riedl, A. (2009). Institution formation in public goods games. *American Economic Review*, 99, 1335-1355.
- Maynard Smith, J. and Szathmari, E. (1995). *The Major Transitions in Evolution*, New York: Oxford University Press.
- Markussen, T., Putterman, L., and Tyran, J-R. (2011). Self-organization for collective action: an experimental study of voting on formal, informal, and no sanction regimes. Working Papers 2011-4, Brown University, Department of Economics.
- Nikiforakis, N. (2008). Punishment and counter-punishment in public good games: can we really govern ourselves? *Journal of Public Economics*, 92, 91-112.
- Nikiforakis, N. and Engelmann, D (2011) Altruistic punishment and the threat of feuds. *Journal of Economic Behavior and Organization*, 78, 319-332.
- Oliver, P. (1980). Rewards and punishments as selective incentives for collective action: theoretical investigations. *American Journal of Sociology*, 85, 1356-1375.
- Olson, M. (1965). *The Logic of Collective Action: Public Goods and the Theory of Groups*, Harvard: Harvard University Press.
- Ostrom, E. (1990). *Governing the Commons: The Evolution of Institutions for Collective Action*, Cambridge: Cambridge University Press.
- Ostrom, E. (2005). *Understanding Institutional Diversity*, Princeton: Princeton University Press.
- Rockenbach, B. and Milinski, M. (2006). The efficient interaction of indirect reciprocity and costly punishment. *Nature*, 444, 718–723.
- Sigmund, K. (2007). Punish or perish? Retaliation and cooperation among humans. *Trends in Ecology and Evolution*, 22, 593-600.
- Sigmund, K. (2010). *The Calculus of Selfishness*, Princeton: Princeton University Press.
- Sigmund, K., De Silva, H., Traulsen, A. and Hauert, C. (2010). Social learning promotes institutions for governing the commons. *Nature*, 466, 861-863.
- Sigmund, K., Hauert, C., Traulsen, A. and De Silva, H. (2011). Social control and the social contract: the emergence of sanctioning systems for collective action. *Dynamic Games and Applications*, 1, 149-171.
- Sugden, R. (1986). *The economics of rights, cooperation and welfare*, Oxford: Blackwell.
- Sutter, M., Haigner, S. and Kocher, M. G. (2010). Choosing the carrot or the stick? Endogenous institutional choice in social dilemma situations. *Review of Economic Studies*, 77, 1540-1566.
- Tiebout, C. (1956). A pure theory of local expenditures. *Journal of Political Economy*, 64, 416-424.

- Traulsen, A., Röhl, T. and Milinski, M. (2012). An economic experiment reveals that humans prefer pool punishment to maintain the commons. *Proceedings of the Royal Society B: Biological Sciences*, 279, 3716-3721.
- Tyler, T. R. and DeGoey, P. (1995). Collective restraint in social dilemmas: procedural justice and social identification effects on support for authorities. *Journal of Personality and Social Psychology*, 69, 482-497.
- Van Vugt, M., Henrich, J. and O'Gorman, R. (2009). Constraining free riding in public good games: designated solitary punishers can sustain human cooperation. *Proceedings of the Royal Society B: Biological Sciences*, 276, 323-329.
- Walker, J. M., Gardner, R., Herr, A. and Ostrom, E. (2000). Collective choice in the commons: experimental results on proposed allocation rules and votes. *The Economic Journal*, 110, 212-34.
- Wilson, D.S. (1975). A theory of group selection. *Proceedings of the National Academy of Sciences*, 72, 13-146.
- Yamagishi, T. (1986). The provision of a sanctioning system as a public good. *Journal of Personality and Social Psychology*, 51, 110-116.

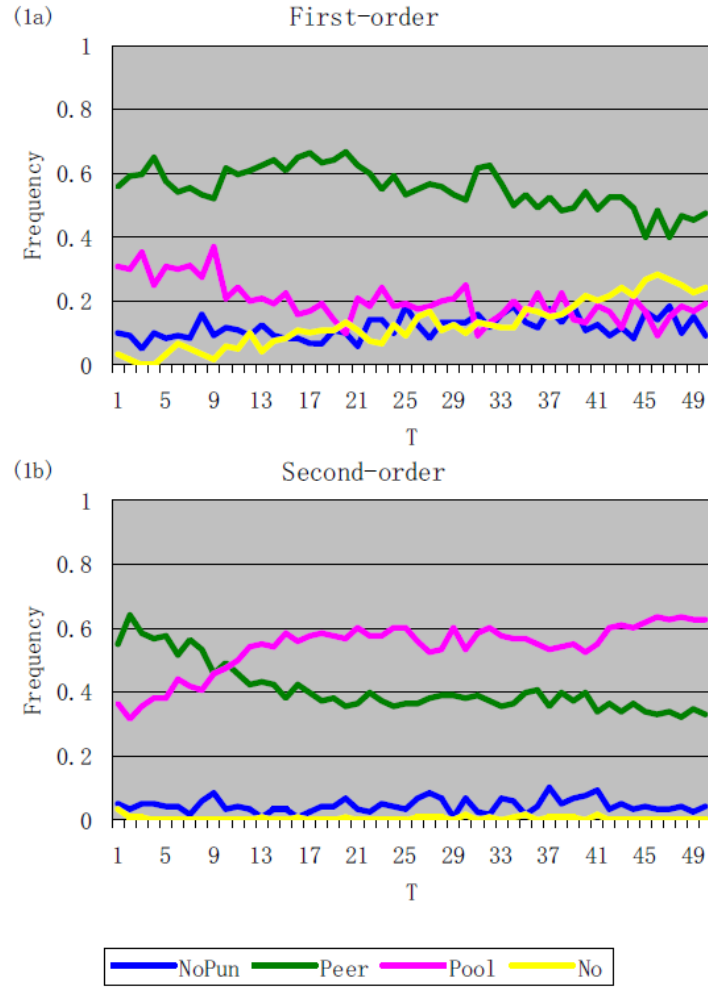


Figure 1

Time-evolution of the frequencies of players voting for the games (NoPun), (Peer), (Pool) or (No) in (a) the first-order treatment and (b) the second order treatment. In each case the numbers are obtained by summing over the nine corresponding groups. (For the frequencies in each group, we refer to Figure S2 of the Supplementary Online Resource).

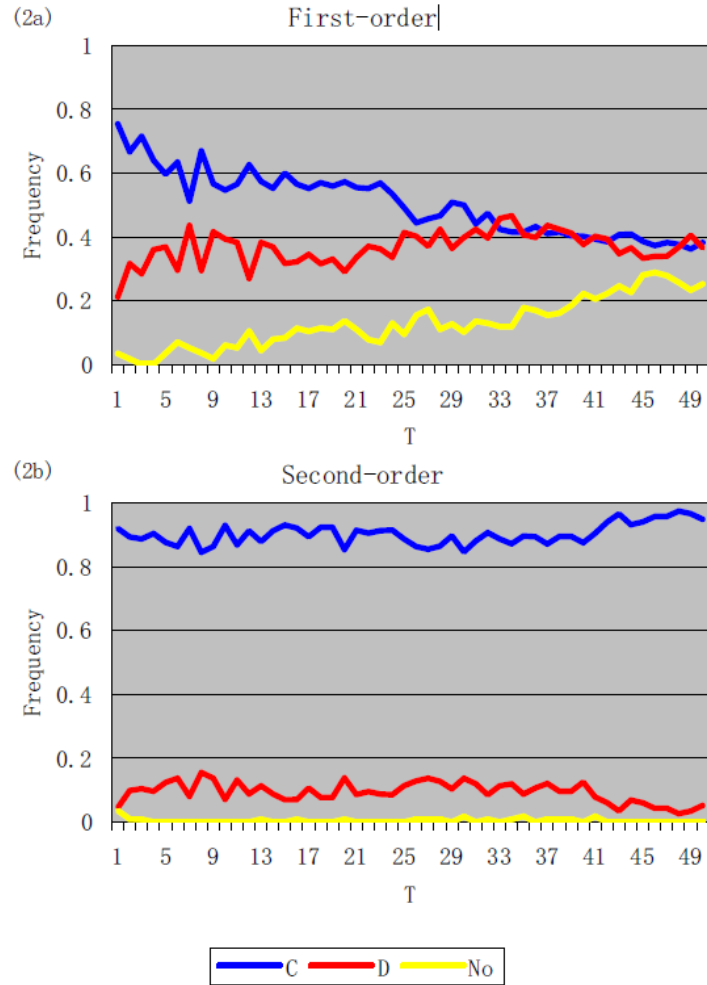


Figure 2

Time-evolution of the frequencies of cooperation (C, blue), defection (D, red) and non-participation (No, yellow) over 50 rounds in the first- and the second-order treatments, summed over the nine corresponding groups. (2a) In the first-order treatment, defection was chosen by about one-third of the players in each round. The number of contributions declined in favor of non-participation. (2b) In the second-order treatment, almost all the players chose to contribute. This cooperative regime was stably sustained.

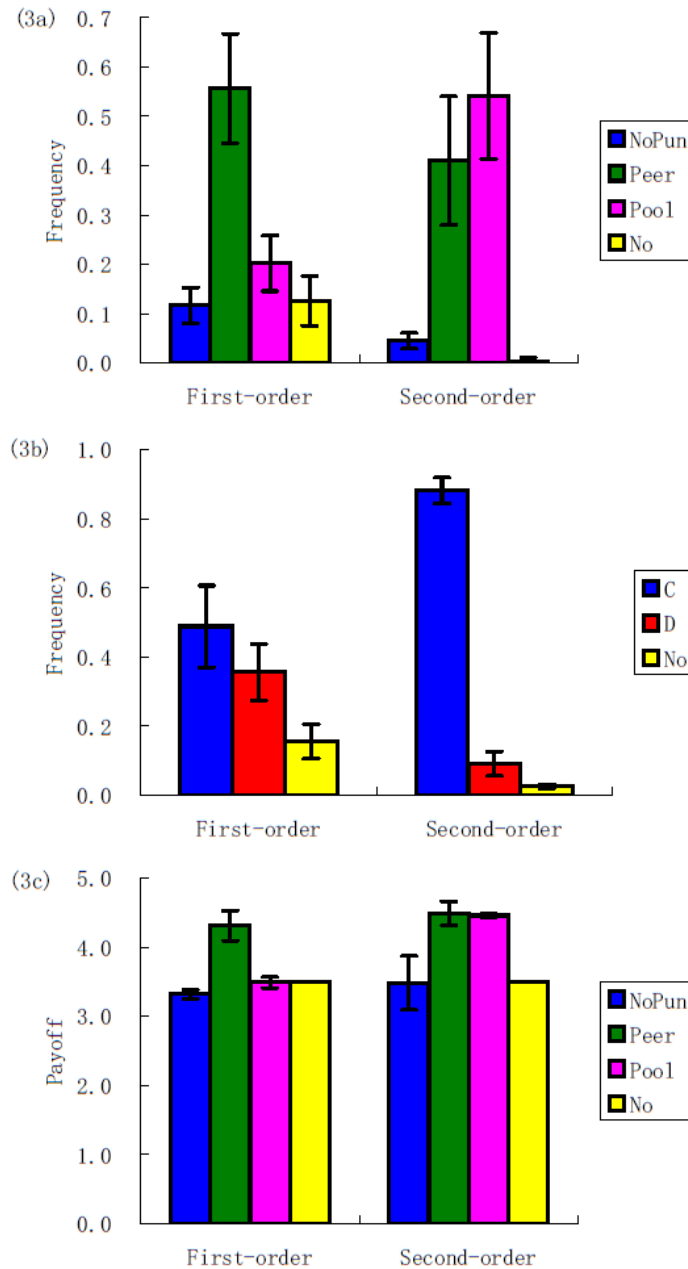


Figure 3

(3a) Frequencies of the decisions in favor of the different games, over 50 rounds, for the first- and the second-order treatments. In the first-order treatment, peer punishment is favored. In the second-order treatment, pool punishment is more frequent, but error bars overlap. (3b) Frequencies of the decisions to contribute to the Mutual Aid game, to defect (i.e., not to contribute) and to opt for non-participation, averaged over 50 rounds. Contribution is strongly promoted in the second-order treatment. (3c) Payoffs obtained for the different games (NoPun), (Peer), (Pool), averaged over fifty rounds, do not greatly differ. Nevertheless, in the first-order treatment, peer punishment games, and in the second-order treatment, both peer and pool punishment games provided the highest average payoff. Error bars indicate standard errors of the group means.