

Democratic decisions establish stable authorities that overcome the paradox of second-order punishment

Christian Hilbe^{a,1}, Arne Traulsen^a, Torsten Röhlf^a, and Manfred Milinski^b

^aEvolutionary Theory Group and ^bDepartment of Evolutionary Ecology, Max Planck Institute for Evolutionary Biology, 24306 Plön, Germany

Edited by Brian Skyrms, University of California, Irvine, CA, and approved November 21, 2013 (received for review August 14, 2013)

Individuals usually punish free riders but refuse to sanction those who cooperate but do not punish. This missing second-order peer punishment is a fundamental problem for the stabilization of cooperation. To solve this problem, most societies today have implemented central authorities that punish free riders and tax evaders alike, such that second-order punishment is fully established. The emergence of such stable authorities from individual decisions, however, creates a new paradox: it seems absurd to expect individuals who do not engage in second-order punishment to strive for an authority that does. Herein, we provide a mathematical model and experimental results from a public goods game where subjects can choose between a community with and without second-order punishment in two different ways. When subjects can migrate continuously to either community, we identify a bias toward institutions that do not punish tax evaders. When subjects have to vote once for all rounds of the game and have to accept the decision of the majority, they prefer a society with second-order punishment. These findings uncover the existence of a democracy premium. The majority-voting rule allows subjects to commit themselves and to implement institutions that eventually lead to a higher welfare for all.

evolution of cooperation | pool punishment | institution formation

The success of collective action and the maintenance of commonly shared infrastructure is often endangered by free riders, subjects who reap the benefits of public goods without contributing to them (1, 2). To mitigate the free riders' destructive potential, many communities install specialized authorities that monitor the subjects' behavior and sanction wrong-doers (3–7). Examples, such as modern courts and the police system, indicate that the maintenance of such institutions is costly. They also constitute a commonly shared infrastructure, which can be exploited just as the original public good that the institution was designed for to protect. Thus, a second-order dilemma arises.

Second-order dilemmas appear in various forms and are considered as a serious obstacle to the evolution of cooperation (8–10). For example, in the absence of a policing authority, group members may take the job onto themselves, punishing others directly. There is overwhelming evidence that subjects are willing to sanction free riders at a cost to themselves (11–14), although individuals typically refuse to exert second-order punishment (15). However, peer punishment can have detrimental consequences on welfare, as the punishment costs may override the benefits of increased cooperation (16) and due to the problems of antisocial punishment (17) and retaliation (18, 19). Peer punishment may pay in the long run, but only when interactions take place in small and stable groups (20). These restrictions may be the reason why modern states have abolished decentralized sanctioning (21).

To explain the transition from decentralized peer punishment to institutional pool punishment (22), recent theoretical and experimental evidence highlights the critical role of second-order punishment (23–26). These studies indicate that such institutions can only persist when they additionally punish individuals who do not support the central authority. The presence of a powerful authority restricts the subjects' strategic options and effectively

forces them to cooperate. As this implies a considerable loss of individual freedom, it is unclear under which conditions subjects would voluntarily submit to such a Leviathan (27). There are different views on this problem: Hardin argued that “we accept compulsory taxes because we recognize that voluntary taxes would favor the conscienceless” (2). However, previous studies have also shown that maintaining costly institutions may result in lower average payoffs (23, 25). Under which conditions would subjects agree to implement a central authority that enforces its continued existence with second-order punishment?

To investigate this question, we conducted an experimental public goods game. The experiment consisted of three independent blocks, each block having several rounds (Table 1 and *Materials and Methods*). During the first two blocks of the experiment, consisting of 10 rounds each, subjects first had to decide whether they want to participate in the game or abstain to secure a small payoff. Participants were then asked whether they want to pay taxes to a central authority and whether they want to cooperate by contributing money to a common pool. If at least one subject paid taxes, the central authority was established and either punished both noncontributors and tax evaders (institution with second-order punishment) or just noncontributors (institution without second-order punishment). If subjects failed to establish such an authority, the public goods game took the form of a conventional social dilemma (mutual cooperation was the optimal outcome for the group, in which case each individual's best choice was to free ride).

In the last block of the experiment, consisting of 15 rounds, subjects had to choose between an authority with or without second-order punishment. As the subjects' choice may depend on the voting mechanism that allows individuals to choose

Significance

Humans usually punish free riders but refuse to sanction those who cooperate but do not punish. However, such second-order punishment is essential to maintain cooperation. The central authorities established in modern societies punish both free riders and tax evaders. This is a paradox: would individuals who do not engage in second-order punishment strive for an authority that does? We address this puzzle with a mathematical model and an economic experiment. When individuals can choose between authorities by migrating between different communities, we find a costly bias against second-order punishment. When subjects use a majority vote instead, they vote for an authority with second-order punishment. These findings also suggest that other pressing social dilemmas could be solved by democratic voting.

Author contributions: C.H., A.T., T.R., and M.M. designed research; C.H. and M.M. performed research; C.H. analyzed data; and C.H., A.T., and M.M. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

Freely available online through the PNAS open access option.

¹To whom correspondence should be addressed. E-mail: hilbe@evolbio.mpg.de.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1315273111/-DCSupplemental.

Table 1. Overview of the experimental design

Treatment	Number of groups	Block I		Block II		Block III
		Rd. 1–5	Rd. 6–10	Rd. 11–15	Rd. 16–20	Rd. 21–35
A	5	Without ZOP	With ZOP	Without ZOP	With ZOP	Foot voting
	5	Without ZOP	With ZOP	With ZOP	Without ZOP	
B	8	Without ZOP	With ZOP	Without ZOP	With ZOP	Majority voting
	7	Without ZOP	With ZOP	With ZOP	Without ZOP	

In the first two blocks of the experiment, subjects gained experience with punishment institutions with and without second-order punishment (ZOP). In the third block, subjects could choose between these two institutional rules. To avoid sequence effects, there are two versions of each treatment. Only the results of blocks II and III are analyzed further. In the subsequent figures, green colors refer to results of block II. Red and blue colors refer to results of block III, for the foot voting treatment and the majority voting treatment, respectively.

between different alternatives (28), we distinguished two different treatments. (A) Subjects can migrate to either community (foot-voting treatment): here, subjects could choose in each round of the last block between an authority with or without second-order punishment. They only interacted with individuals who chose the same institutional rule. Previous experiments used such a voting scheme to show that humans prefer peer punishment institutions to punishment-free institutions (13), even if reputation allows for an alternative mechanism to govern the commons (14). (B) Subjects participate in a democratic vote (majority-voting treatment): subjects had to vote for their preferred institution in the beginning of the last block. The institutional rule that obtained a majority of votes was then implemented for all remaining 15 rounds and was imposed on all group members. Such a scheme of elected authorities can elicit higher contributions to public goods than randomly chosen authorities (29) and help subjects to coordinate on pool punishment systems with optimal parameters (30).

Results

Based on a theoretical model (described in detail in the *SI Text*), we expected that only a minority of subjects would pay taxes if there is no punishment for tax evasion. As a consequence, we also predicted that authorities without second-order punishment would result in less cooperation and lower average payoffs. An analysis of block II of our experiments (in which subjects could not choose between different institutions) confirms these predictions (Fig. 1). Second-order punishment institutions facilitated higher average payoffs (payoffs increased from 0.63 to 0.86 Euro per round when tax evaders were punished, Wilcoxon matched-pairs signed-rank test, $n_{A+B} = 25$, $Z = 4.023$, $P < 0.001$; we used two-tailed statistics throughout), and led to more cooperation (the fraction of cooperators increased from 61.4% to 96.5%, Wilcoxon matched-pairs signed-rank test, $Z = 4.270$, $P < 0.001$). This efficiency advantage of second-order punishment institutions suggests that subjects should prefer this institutional rule when given a choice in the last block, independent of the implemented voting rule.

However, for the votes before the first round of the last block, we observed a significant treatment effect (Kolmogorov-Smirnov two-sample test, $n_A = 10$, $n_B = 15$, $K = 1.470$, $P = 0.027$). In foot-voting groups, subjects initially preferred institutions without second-order punishment (with 8 of 10 groups having a majority against second-order punishment in the first round of block III; Fig. 2A). Only the groups in the majority-voting treatment showed a clear preference for second-order punishment, with 12 of 15 groups voting for the respective institution (binomial test, $n_B = 15$, $P = 0.035$; Fig. 2B). Over the course of the experiment, this treatment effect waned; by the end of the last block, in four more groups in the foot-voting treatment, the majority of players switched to second-order punishment. In total, 35.6% of the subjects in the foot-voting treatment played under an authority

with second-order punishment compared with the 80.0% in the majority-voting treatment (Fig. 2C).

To study the dynamics during the third block, we compared the players' strategies in the beginning (rounds 1–5) with the strategies in the end (rounds 11–15) (Fig. 3). Although behaviors in the majority-voting treatment were stable as predicted (none of the considered variables changed significantly over time), we found significant learning effects in the foot-voting treatment. Driven by a stronger preference for second-order punishment (rounds 1–5, 19.6%; rounds 11–15, 54.8%; Wilcoxon matched-pairs signed-rank test, $Z = 2.429$, $P = 0.015$), we found a significant increase in the number of tax payers (rounds 1–5, 18.4%; rounds 11–15, 55.6%; Wilcoxon matched-pairs signed-rank test, $Z = 2.374$, $P = 0.018$). The higher willingness to pay taxes resulted in a reduction of the number of defectors (rounds 1–5, 27.2%; rounds 11–15, 13.2%; Wilcoxon matched-pairs signed-rank test, $Z = -2.095$, $P = 0.036$), whereas it had no significant impact on the number of cooperators or on the resulting average payoff.

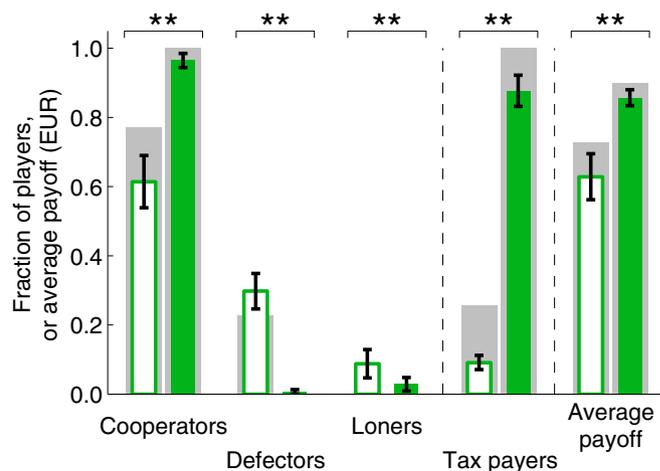


Fig. 1. Second-order punishment promotes cooperation. The graph shows the fraction of cooperators (individuals who contribute to the common pool), defectors (individuals who do not contribute), loners (individuals who decide not to participate in the game), and the fraction of subjects paying taxes, as well as the resulting average payoff. Colored bars depict the experimental results of block II, with empty bars corresponding to rounds without second-order punishment, and filled bars showing rounds with second-order punishment. Two stars indicate significance at the $\alpha = 0.01$ level (using Wilcoxon matched-pairs signed-rank tests). Gray bars depict the theoretical predictions based on a social learning model for the one-shot game (see *SI Text* for details). As predicted, second-order punishment resulted in more cooperation in the public goods game, as more subjects were willing to support the punishment institution by paying taxes. Overall, this led to a significant increase of average payoffs.

only the subjects in the majority-voting treatment succeeded in implementing the corresponding authority. Subjects in the foot-voting treatment showed a costly bias in favor of institutions without second-order punishment. Various causes may be responsible for this positive effect of the majority-voting rule. First, the outcome of the democratic decision was binding for all 15 rounds of the last block; this may have triggered subjects to take a long-run perspective and to anticipate the risks and benefits of each punishment institution. Second, the decision to migrate to either community in the foot-voting treatment only affects each subject individually. When using a majority vote instead, individuals can bind each other. This option to bind each other may have triggered subjects to take a group perspective and to opt for the institution that leads to a beneficial group dynamics, rather than choosing an institution that promises individual advantages. Buchanan and Congleton argued that “persons agree to constraints on their own liberties in exchange for comparable constraints being imposed on the liberties of others.” (35) The majority-voting rule can be seen as a mechanism that helps individuals to implement such beneficial constraints.

Institutions are inherently unstable when they only apply to a subset of community members (36). A similar problem arises if institutions are only funded by such a subset (23–26, 37): when paying taxes occurs on a voluntary basis, tax evasion can lead to the breakdown of cooperation (as also shown in Fig. 1). Thus, the stability of many modern institutions requires second-order punishment, where subjects that do not support the central institution are punished just as ordinary free riders. Interestingly, the delegation of punishment to central institutions may in turn facilitate the emergence of second-order punishment: First, since institutions need to be funded in advance, second-order free riders (i.e., tax evaders) are easy to detect (23). Second, setting up a punishment institution to protect a community from wrongdoers may be expensive, but once the institution is established, extending its scope to prosecute also tax evaders seems to be relatively cheap. In this way, first-order and second-order punishment become linked: the same institution automatically engages in both forms of punishment. This linkage is critical for the maintenance of second-order punishment, as it removes the need for further levels of punishment (such as third-order punishment) to stabilize the lower levels (38). In peer punishment, it is not immediately clear how such a linkage between first-order and second-order punishment could evolve (33, 34). When punishment is delegated to a central authority instead, this linkage can be implemented easily.

We showed here that a pool punishment regime with second-order punishment can emerge if individuals have the freedom to bind each other with a majority vote, but not if they can individually reconsider their decision after each round. In our experiments, democracy prompts individuals to commit them-

selves and to make institutional choices that enhance the welfare of all.

Materials and Methods

Experimental Design. Experiments were conducted in November 2012 at the University of Kiel, Kiel, Germany, with 125 subjects recruited from a first-year course in biology. Twenty-five groups of five subjects played 35 rounds of a public goods game. In each round, players first had to choose between being a loner (fixed payoff of € 0.40) and taking part in the public goods game. Those subjects who decided to participate were then asked whether they want to pay taxes for a punishment institution. Individual taxes depended on the number of tax payers (which is a typical feature of models of coordinated punishment) (39, 40): if i is the number of tax payers, then the tax was set to $0.05 + 0.45/i$ (institution without second-order punishment) and to $0.05 + 0.5/i$ (institution with second-order punishment), respectively. These parameters reflect our assumption that a punishment institution comes with high fixed costs, but comparably low variable costs (i.e., extending the institution’s scope to punish also tax evaders does not duplicate the costs of the institution). If at least one participant paid taxes, the punishment institution was established and either imposed a fine on defectors only (institution without second-order punishment) or on defectors and tax evaders (institution with second-order punishment). The fine for defectors (and tax evaders) was set to € 1.00 for each offense. Subjects were informed about whether someone paid taxes before they had to decide whether they want to contribute € 0.50 to a common pool. Total contributions to the pool were multiplied by 3.1 and redistributed to all group participants. See *SI Text* for further details.

Theoretical Predictions. To illustrate the possible strategic considerations of the players, let us calculate the symmetric subgame perfect equilibria for the one-shot public good game. (i) Without second-order punishment, the decision to pay taxes becomes a volunteer’s dilemma (41–43): subjects benefit from the presence of a punishment authority, but they want others to pay the taxes. The symmetric solution to this dilemma is to pay taxes with a certain probability q_T . This probability can be calculated by comparing the expected cost of paying taxes with the expected loss to be in a group where no one pays taxes (and hence no one cooperates)

$$\sum_{i=0}^4 \binom{4}{i} q_T^i (1 - q_T)^{4-i} \cdot \left(0.05 + \frac{0.45}{1+i}\right) = 1.05 \cdot (1 - q_T)^4. \quad [1]$$

Solving this equation leads to the prediction that all players participate in the game, pay taxes with probability $q_T = 25.6\%$, and contribute in case there was at least someone who paid taxes. In this equilibrium, players earn on average € 0.73 per round. (ii) With second-order punishment, payoff dominance suggests that players participate in the game, pay taxes, and contribute to the common pool. This equilibrium results in an expected payoff of € 0.90. Therefore, independent of the voting procedure, equilibrium payoffs are higher under second-order punishment. These static predictions are also confirmed by a dynamic learning model (*SI Text*).

ACKNOWLEDGMENTS. We thank M. Abou Chakra and two anonymous referees for insightful comments. We thank K. Hagel and H. Brendelberger for support performing the experiment and the 125 students for their participation.

- Olson M (1971) *The Logic of Collective Action: Public Goods and the Theory of Groups* (Harvard Univ Press, Cambridge, MA).
- Hardin G (1968) The tragedy of the commons. *Science* 162(3859):1243–1248.
- Yamagishi T (1986) The provision of a sanctioning system as a public good. *J Pers Soc Psychol* 51(1):110–116.
- Ostrom E (1990) *Governing the Commons: The Evolution of Institutions for Collective Action* (Cambridge Univ. Press, Cambridge, UK).
- O’Gorman R, Henrich J, Van Vugt M (2009) Constraining free riding in public goods games: designated solitary punishers can sustain human cooperation. *Proc Biol Sci* 276(1655):323–329.
- Pinker S (2011) *The Better Angels of Our Nature: Why Violence Has Declined* (Penguin Books, New York).
- Sasaki T, Brännström Å, Dieckmann U, Sigmund K (2012) The take-it-or-leave-it option allows small penalties to overcome social dilemmas. *Proc Natl Acad Sci USA* 109(4):1165–1169.
- Fowler JH (2005) Human cooperation: Second-order free-riding problem solved? *Nature* 437(7058):E8.
- Sigmund K (2007) Punish or perish? Retaliation and collaboration among humans. *Trends Ecol Evol* 22(11):593–600.
- Hilbe C, Traulsen A (2012) Emergence of responsible sanctions without second order free riders, antisocial punishment or spite. *Sci Rep* 2:458.
- Fehr E, Gächter S (2002) Altruistic punishment in humans. *Nature* 415(6868):137–140.
- Henrich J, et al. (2006) Costly punishment across human societies. *Science* 312(5781):1767–1770.
- Gürerk Ö, Irlenbusch B, Rockenbach B (2006) The competitive advantage of sanctioning institutions. *Science* 312(5770):108–111.
- Rockenbach B, Milinski M (2006) The efficient interaction of indirect reciprocity and costly punishment. *Nature* 444(7120):718–723.
- Kiyonari T, Barclay P (2008) Free-riding may be thwarted by second-order rewards rather than punishment. *J Pers Soc Psychol* 95:826–842.
- Dreber A, Rand DG, Fudenberg D, Nowak MA (2008) Winners don’t punish. *Nature* 452(7185):348–351.
- Herrmann B, Thöni C, Gächter S (2008) Antisocial punishment across societies. *Science* 319(5868):1362–1367.
- Nikiforakis N (2008) Punishment and counter-punishment in public good games: Can we really govern ourselves? *J Public Econ* 92(1-2):91–112.
- Fehl K, Sommerfeld RD, Semmann D, Krambeck HJ, Milinski M (2012) I dare you to punish me—vendettas in games of cooperation. *PLoS ONE* 7(9):e45093.
- Gächter S, Renner E, Sefton M (2008) The long-run benefits of punishment. *Science* 322(5907):1510.
- Guala F (2012) Reciprocity: Weak or strong? What punishment experiments do (and do not) demonstrate. *Behav Brain Sci* 35(1):1–15.

