

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

Christian Hilbe<sup>a,1</sup>, Arne Traulsen<sup>a</sup>, Torsten Röhlf<sup>a</sup>, and Manfred Milinski<sup>b</sup>

<sup>a</sup>Evolutionary Theory Group and <sup>b</sup>Department 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.

<sup>1</sup>To 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](http://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).

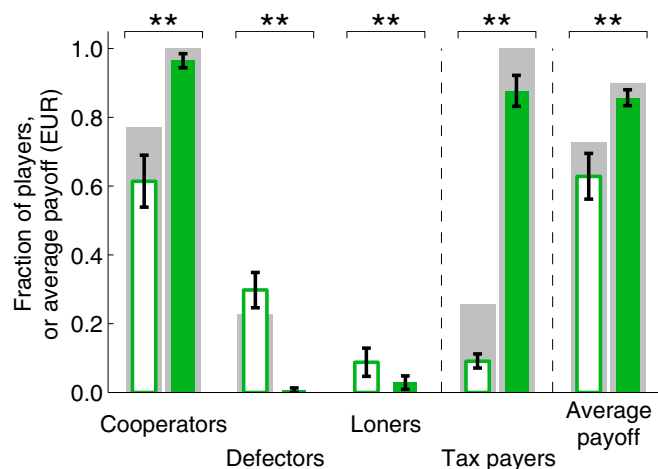
**Results**

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.

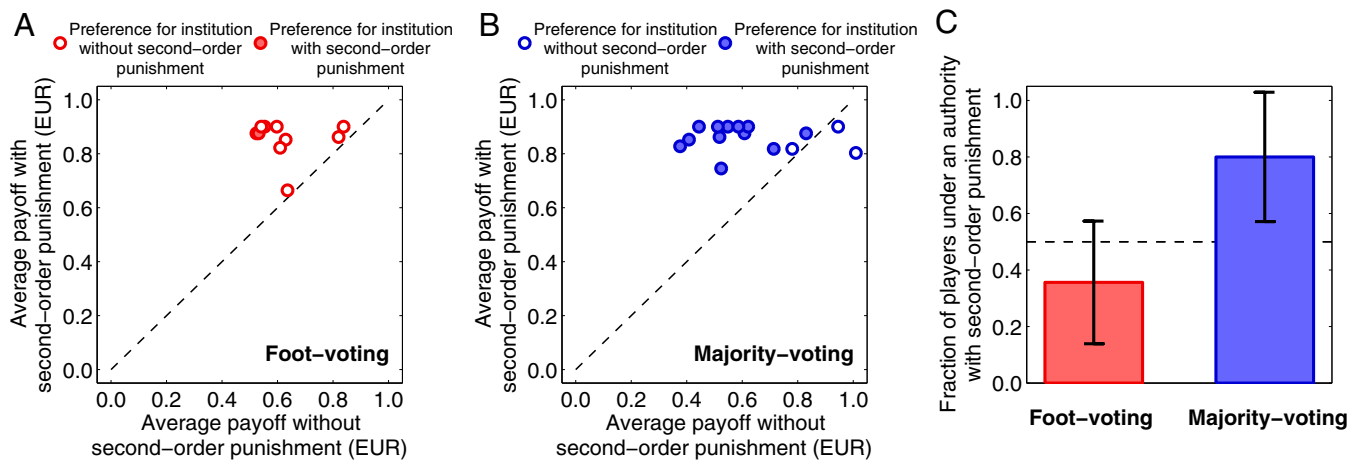
However, for the votes before the first round of the last block, we observed a significant treatment effect (Kolmogorov-Smirnov two-sample test,  $n_A = 10$ ,  $n_B = 15$ ,  $K = 1.470$ ,  $P = 0.027$ ). In foot-voting groups, subjects initially preferred institutions without second-order punishment (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.



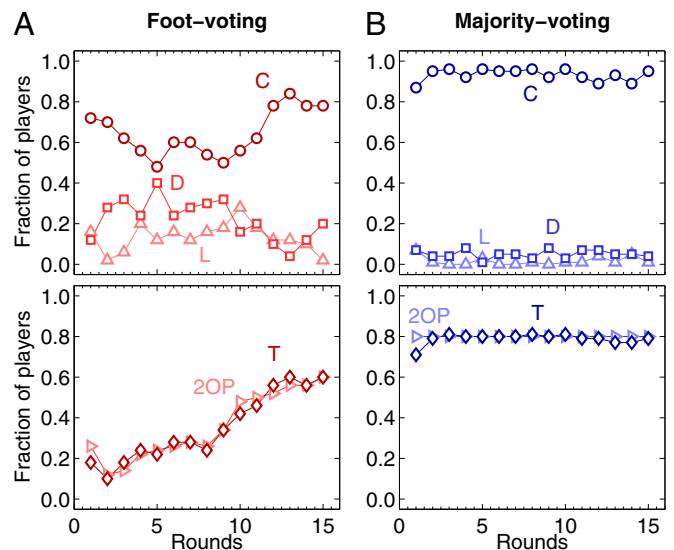
**Fig. 2.** Democratic decisions establish institutions with second-order punishment. (A and B) Horizontal axis shows the average payoff that each group obtained during block II when acting under an institution without second-order punishment, whereas the vertical axis shows the average payoff obtained during rounds with second-order punishment. Thus, for all groups above the diagonal, second-order punishment resulted in higher payoffs during block II. The filling of the symbol indicates the group's vote in the first round of block III: the group symbol is filled if at least half of the group members voted for second-order punishment. Basing decisions on efficiency would imply that symbols above the dashed line are filled. (A) In the foot-voting treatment, only 2 of 10 groups had at least half of the players voting for an institution with second-order punishment (although all 10 groups earned, on average, more under such an institution). (B) Under the majority-voting rule, 12 of 15 groups voted for second-order punishment. (C) Over all 15 rounds of block III, subjects in the foot-voting treatment chose second-order punishment in 35.6% of all cases. In contrast, 80% of the subjects in the majority-voting treatment were governed by an institution with second-order punishment (see *SI Text* for further details on the subject's voting behavior).

Overall, payoffs in the foot-voting treatment were clearly below the payoffs under majority voting (€ 0.64 compared with € 0.86; Mann-Whitney *U* test,  $n_A = 10$ ,  $n_B = 15$ ,  $Z = 3.534$ ,  $P < 0.001$ ). These findings indicate that foot voting has led to a costly bias in favor of unstable institutions without second-order punishment. This bias decreased over time, but it did not disappear completely. The lower payoffs in the foot-voting treatment are unexpected, as this treatment can give rise to smaller group sizes (subjects only interacted with those who voted for the same institution), which in turn facilitates cooperation (31). Post-experiment questionnaires suggest that subjects in the foot-voting treatment hoped that eventually everyone would play fair, such that there would be no need for a punishment institution. Typically, this hope did not come true; none of the groups in the foot-voting treatment and only one group in the majority-voting treatment managed to choose an institution without second-order punishment and to cooperate in all 15 rounds without paying taxes. The majority-voting rule, on the other hand, seemed to trigger subjects to make group-beneficial decisions. Indeed, for 14 of the 15 groups, the majority vote resulted in the institutional rule that proved to be more efficient in the preceding blocks (Fig. 2B; binomial test,  $P < 0.001$ ).

### Discussion

Economic experiments have repeatedly shown that individuals are willing to punish free riders (11–14). At the same time subjects either abstain from second-order punishment opportunities (15), or they may even abuse them for perverse punishment (i.e., forms of punishment that have the effect of reducing future cooperation) (32). The subjects' reluctance to sanction nonpunishers is surprising because second-order punishment is of fundamental importance for the stability of decentralized peer punishment (33, 34), and it also plays a crucial role for the evolution of central pool punishment institutions (23–26). Indeed, in most societies today, central punishment institutions are funded by compulsory taxes rather than by voluntary contributions. The emergence of such authorities poses a puzzle: does this mean that individuals are able to implement institutions with second-order punishment, although the individuals themselves are not willing to engage in second-order punishment?

With an economic experiment, we investigated two possible routes for the emergence of central institutions with second-order punishment. In one treatment, subjects could choose between different institutions by migrating to the respective community, whereas in the other treatment, subjects could vote for their preferred institution in a democratic election. Our mathematical model of social learning (see *SI Text* for details) suggests in both cases that institutions with second-order punishment are more stable and result in higher payoffs. However, in our experiment



**Fig. 3.** Social learning leads to a trend toward second-order punishment in block III. (Upper) Fraction of cooperators C, defectors D, and loners L in each round. (Lower) Fraction of players preferring an institution with second-order punishment (2OP), as well as the fraction of players who paid taxes (T). (A) Foot-voting groups learned to adopt institutions with second-order punishment, which led to an increase in the number of taxpayers, and to a reduction of the number of defectors. (B) In the majority-voting treatment, on the other hand, behaviors were stable.

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.

22. Witt U (1992) The emergence of a protective agency and the constitutional dilemma. *Constitutional Political Economy* 3:255–266.
23. Sigmund K, De Silva H, Traulsen A, Hauert C (2010) Social learning promotes institutions for governing the commons. *Nature* 466(7308):861–863.
24. Perc M (2012) Sustainable institutionalized punishment requires elimination of second-order free-riders. *Sci Rep* 2:344.
25. Traulsen A, Röhl T, Milinski M (2012) An economic experiment reveals that humans prefer pool punishment to maintain the commons. *Proc Biol Sci* 279(1743):3716–3721.
26. Zhang B, Li C, De Silva H, Bednarik P, Sigmund K (2013) The evolution of sanctioning institutions: An experimental approach to the social contract. *Exper Econ*, 10.1007/s10683-013-9367-7.
27. Hobbes T (1651) *Leviathan* (Andrew Crooke, London).
28. Zeckhauser R (1973) Voting systems, honest preferences and Pareto optimality. *Am Polit Sci Rev* 67(3):934–946.
29. Baldassarri D, Grossman G (2011) Centralized sanctioning and legitimate authority promote cooperation in humans. *Proc Natl Acad Sci USA* 108(27):11023–11027.
30. Putterman L, Tyran J, Kamei K (2011) Public goods and voting on formal sanction schemes. *J Public Econ* 95(9–10):1213–1222.
31. Ledyard JO (1995) *The Handbook of Experimental Economics*, eds Kagel JH, Roth AE (Princeton Univ. Press, Princeton, NJ).
32. Cinyabuguma M, Page T, Putterman L (2006) Can second-order punishment deter perverse punishment? *Exp Econ* 9(3):265–279.
33. Axelrod R (1986) An evolutionary approach to norms. *Am Polit Sci Rev* 80(4):1095–1111.
34. Colman AM (2006) The puzzle of cooperation. *Nature* 440(7085):744–745.
35. Buchanan JM, Congleton RD (1998) *Politics by Principle, Not Interest* (Cambridge Univ Press, Cambridge, UK).
36. Kosfeld M, Okada A, Riedl A (2009) Institution formation in public goods games. *Am Econ Rev* 99(4):1335–1355.
37. Schoenmakers S (2013) Pool-punishment and opportunistic cooperation in voluntary and compulsory games. An evolutionary game theory model. MS thesis (Univ of Oldenburg, Oldenburg, Germany).
38. Yamagishi T, Takahashi N (1994) *Social Dilemmas and Cooperation*, eds Schulz U, Albers W, Mueller U (Springer, Berlin), pp 311–326.
39. Boyd R, Gintis H, Bowles S (2010) Coordinated punishment of defectors sustains cooperation and can proliferate when rare. *Science* 328(5978):617–620.
40. Dercole F, De Carli M, Della Rossa F, Papadopoulos AV (2013) Overpunishing is not necessary to fix cooperation in voluntary public goods games. *J Theor Biol* 326:70–81.
41. Diekmann A (1985) Volunteer's dilemma. *J Conflict Resolut* 29(4):605–610.
42. Raihani NJ, Bshary R (2011) The evolution of punishment in n-player public goods games: a volunteer's dilemma. *Evolution* 65(10):2725–2728.
43. Przepiorka W, Diekmann A (2013) Individual heterogeneity and costly punishment: A volunteer's dilemma. *Proc Biol Sci* 280(1759):20130247.

# Supporting Information

Hilbe et al. 10.1073/pnas.1315273111

## SI Text

**Experimental Methods.** The experiment was conducted in November 2012. In total, we recruited 125 volunteers from a first-year course in biology at the University of Kiel (Kiel, Germany). These volunteers participated in 25 groups of five subjects each in a computerized experiment. Subjects were separated by opaque partitions. They received the instructions, and they communicated their decisions with a laptop computer. Before each experimental session, subjects were orally informed about how to operate the computers and about the measures that were taken to ensure their absolute anonymity (which included that participants made their decisions under a neutral pseudonym and that they were not allowed to talk to each other during or after the experiment). Moreover, they were informed that they played for real money and that the payment of their earnings was organized in a way that fully maintains the subjects' anonymity (1).

The experimental game used a minimalist model of a public goods game with punishment institutions. In the first two blocks, individuals could decide whether or not to participate in the game, whether or not to pay taxes for the punishment institution, and whether or not to contribute to the common pool. In both blocks, subjects played five rounds under an institution without second-order punishment and five rounds under an institution with second-order punishment. These first two blocks of the experiment were meant to familiarize the subjects with the game and with the consequences of their decisions. In the third block, individuals had to make the additional decision whether they prefer an institution with or without second-order punishment. To avoid sequence effects that may affect behaviors in the last block of the experiment, we considered two different versions of each treatment (in which the game was played either with or without second-order punishment before block III; Table 1).

The subjects were unaware of the number of rounds in each block, but they were informed that the money that they had gained in the previous blocks was safe in their account. The initial endowment of each player was chosen such that each player's overall payoff would sum up to a positive amount after any possible history of the game; thus, players received 12 Euros in the beginning of each of the first two blocks and 18 Euros in the beginning of the last block. Each experiment took ~1.5 h, and subjects earned on average 68.59 Euros. In total, both treatments consisted of 35 (10 + 10 + 15) rounds of the public goods game (Table 1). Each block was introduced by a series of text pages explaining the rules and giving examples of possible outcomes. After each text page, subjects had to confirm that they had read all instructions by clicking on the respective button. A detailed description of the experimental procedure of each treatment is given below:

*i*) Treatment A (foot voting): In each round of block I and block II, players first had to choose between being a loner (fixed payoff of € 0.40) and taking part in the actual public goods game. Those subjects who decided to participate were then asked whether they want to pay taxes for the punishment institution. The total costs of the punishment institution were fixed to € 0.45 (institution without second-order punishment) and € 0.50 (institution with second-order punishment), respectively. As a consequence, the individual tax decreased in the number of tax payers  $i$ ; it was set to  $0.05 + 0.45/i$  (without second-order punishment) and  $0.05 + 0.5/i$  (with second-order punishment), respectively.

The constant term 0.05 reflects the administration costs per tax payer. After the decision to pay taxes, players were informed whether there was at least one player who paid for the punishment institution. Thereafter, players had to decide whether or not to contribute € 0.50 to the common pool, knowing that the total money in this pool is multiplied by 3.1 and then equally shared among all group participants. If at least one player paid taxes, the punishment institution was established, such that defectors (and tax evaders in rounds with second-order punishment) had to pay a fine. The fine was set to € 1.00 for each offense (as a consequence, when a second-order punishment institution was established, a player who neither paid taxes nor contributed to the public pool was punished twice). By the end of each round, subjects were informed about all co-player's decisions and about the resulting payoffs.

In each round of block III, participants first had to decide whether they want to play under an institution with or without second-order punishment. In the subsequent public goods game (which had the same rules as described above), individuals only interacted with subjects who voted for the same punishment institution. All decisions were made by clicking on the respective button on the computer screen. If a player happened to be alone in the respective group (either because the other players decided to be loners, or because they voted for the other institution), then this player automatically became a loner with fixed payoff € 0.40.

*ii*) Treatment B (majority voting): The rules for block I and block II are the same as in the first treatment. After block II, one of the experimenters (M.M.) announced that there would be an election before the last block of the game. In this election, players would be able to vote between an institution with or without second-order punishment, and the institution with the majority of votes would be implemented for the whole group during the entire last block of the experiment. For the election, all subjects had two pieces of paper on their desk, one with the text "Fine for noncontributors" and one with the text "Fine for noncontributors and fine for tax evaders." Subjects voted by putting one of the two pieces of paper into a ballot box (it was ensured that neither the other players nor the experimenters were informed about the individual votes during or after the experiment). The institutional rule that was chosen by a majority of subjects was then implemented for all remaining rounds of block III.

In summary, the treatments only differed in the applied voting rule during block III. Thus, a comparison of these two treatments allows us to investigate the impact of different voting rules on the emergence of institutions with second-order punishment. To study the role of efficiency and of the voting procedure, we analyzed the subjects' behavior in the second block and in the third block.\* We considered groups of subjects as statistical units and used two-tailed tests throughout.

\*The results in block I are qualitatively similar to the results in block II. In particular, aggregating the results of block I and block II would not alter the conclusions of the main text.

**Theoretical Model.** To obtain an intuitive understanding of the possible strategic considerations of the subjects, we first perform a static analysis of the one-shot public goods game with and without second-order punishment, based on the subgame perfect equilibrium concept (2, 3). However, as it has been demonstrated that evolutionary learning processes do not necessarily settle at the subgame perfect equilibrium (4), we complement this static analysis with a social learning model (5, 6). We show that in the limit of strong selection, both models give the same prediction for our experiment. To this end, we abstract from the specific payoff values used in the experiment, and introduce general parameters as explained in Table S1.

**Static analysis.** In the following, we calculate the symmetric subgame perfect equilibrium for the one-shot public goods game with an institution with and without second-order punishment.

i) Without second-order punishment, paying taxes becomes a volunteer's dilemma: if no player pays taxes, the punishment institution will not be established, and thus rational coplayers will not cooperate (resulting in a payoff of zero). On the other hand, if at least one of the players pays taxes, it becomes optimal to cooperate for all players (resulting in a payoff of  $\pi = rc - c - \text{taxes} = 1.05 - \text{taxes}$ ). The symmetric solution to this volunteer's dilemma is to pay taxes with a certain probability  $q_T$  (7). The probability  $q_T$  can be calculated by comparing the expected costs of paying taxes with the expected loss to be in a group of  $n - 1 = 4$  other players where no one pays taxes (and hence no one cooperates)

$$\sum_{i=0}^{n-1} \binom{n-1}{i} q_T^i (1-q_T)^{n-1-i} \cdot \left( \gamma_0 + \frac{\gamma_1}{1+i} \right) = (rc - c) \cdot (1-q_T)^{n-1}. \quad [S1]$$

Solving this equation for the parameters of the experiment yields 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, the probability to cooperate in a given round is  $q_C = 1 - (1-q_T)^n = 77.2\%$ , yielding an expected payoff per round of

$$\pi = (rc - c) \cdot \left( 1 - (1-q_T)^n \right) = 0.73 \text{ Euros}^\dagger \quad [S2]$$

ii) With second-order punishment, the decision to pay taxes is no longer a volunteer's dilemma. Instead, using the symmetric subgame equilibrium, we predict that all players participate, pay taxes, and contribute to the common pool. The average payoff per round becomes  $\bar{\pi} = rc - c - \text{taxes} = 1.05 - .05 - 0.5/5 = 0.90$  Euros.

In equilibrium, players thus earn higher payoffs if they implement an institution with second-order punishment. This conclusion remains unchanged even if we assume that an institution with second-order punishment is twice as costly as an institution without. In that case, the equilib-

rium payoff with second-order punishment becomes  $\bar{\pi} = 1.05 - .05 - 0.9/5 = 0.82$ , which is still above the expected payoff for an institution without second-order punishment,  $\pi = 0.73$ .<sup>‡</sup>

**Social learning model.**

**Setup of the game.** Our social learning model is based on a straightforward application of evolutionary game theory in finite populations (8–10) and comparable to previous models for the evolution of peer and pool punishment (6, 11, 12).

During block I and block II of the experiment, players can make the following decisions:

- i) They can either participate in the public goods game or they can abstain from it;
- ii) Participants are then asked whether they want to pay taxes; and
- iii) Depending on whether someone paid taxes for the punishment institution, participants then have to decide whether to contribute to the common pool.

Overall, such a game allows for seven different strategies, as summarized in Table S2. The strategy  $S_0$  corresponds to a loner, whereas the strategy  $S_1$  may be considered as a selfish player (neither paying taxes, nor contributing to public good). The strategy  $S_2$  is opportunistic, by only contributing to the common pool if a punishment institution has been established. In contrast, a (somewhat paradoxical)  $S_3$  player only contributes if no punishment institution has been established, and an  $S_4$  player always contributes to the common pool but never pays taxes. A player with strategy  $S_5$  can be considered as a righteous citizen, contributing to the common pool and paying taxes for the punishment institution. Last, an  $S_6$  player behaves rather counterintuitively by paying taxes for a punishment institution but not contributing to the common pool.

**Derivation of the payoffs for institutions without second-order punishment.** Let us assume that the game is played in a population of size  $M$  and that the number of players that apply strategy  $S_i$  is given by  $M_i$ , such that  $M_0 + \dots + M_6 = M$ . From this population,  $n \leq M$  players are randomly sampled to play the public goods game. To calculate the expected payoff of each group member, let us assume that the focal player is in a group with  $n_i$  players with strategy  $i$ , such that  $n_0 + \dots + n_6 = n - 1$ . Let  $\vec{n} = (n_0, \dots, n_6)$  denote the vector of these numbers. We need to distinguish three different cases:

- i) All other group members are loners; in this case the player's payoff is  $\sigma$ .
- ii) There is at least one other participant among the other group members, but none of the group members pays taxes; in this case there will be no policing institution, and only players with strategies  $S_3$  and  $S_4$  will contribute to the common pool.
- iii) There is at least one other participant among the other group members and at least one participant who pays taxes; in this case a policing institution is established, and all players with strategies  $S_2, S_4,$  and  $S_5$  contribute to the common pool, whereas the remaining participants with strategies  $S_1, S_3,$  and  $S_6$  are punished.

Considering these three possible cases, expected payoffs  $\pi_i$  are given by  $\pi_0 = \sigma$  and

<sup>†</sup>Additionally to this outcome, the game has a second symmetric subgame perfect equilibrium where all players decide not to participate in the game (yielding the loner's payoff € 0.40). This second equilibrium, however, is dominated in the sense that against rational coplayers, being a loner is a weakly dominated strategy. The same applies to the game with second-order punishment; we will thus neglect this second subgame perfect equilibrium in the following.

<sup>‡</sup>This prediction is based on equilibrium payoffs. Out of equilibrium, when only a small share of players pays taxes, increasing the costs of second-order punishment may have a strong demotivating effect on tax payers. Thus, the incentives to vote for second-order punishment may change more drastically if players are unable to coordinate on the payoff-dominant equilibrium.

$$\begin{aligned}
\pi_1 &= p_1(L)\sigma + \sum_{\substack{n_0 < n-1 \\ n_5+n_6=0}} p_1(\vec{n}) \frac{n_3+n_4}{n-n_0} rc \\
&+ \sum_{n_5+n_6>0} p_1(\vec{n}) \left( \frac{n_2+n_4+n_5}{n-n_0} rc - \beta \right) \\
\pi_2 &= p_2(L)\sigma + \sum_{\substack{n_0 < n-1 \\ n_5+n_6=0}} p_2(\vec{n}) \frac{n_3+n_4}{n-n_0} rc \\
&+ \sum_{n_5+n_6>0} p_2(\vec{n}) \left( \frac{n_2+n_4+n_5+1}{n-n_0} rc - c \right) \\
\pi_3 &= p_3(L)\sigma + \sum_{\substack{n_0 < n-1 \\ n_5+n_6=0}} p_3(\vec{n}) \left( \frac{n_3+n_4+1}{n-n_0} rc - c \right) \\
&+ \sum_{n_5+n_6>0} p_3(\vec{n}) \left( \frac{n_2+n_4+n_5}{n-n_0} rc - \beta \right) \\
\pi_4 &= p_4(L)\sigma + \sum_{\substack{n_0 < n-1 \\ n_5+n_6=0}} p_4(\vec{n}) \left( \frac{n_3+n_4+1}{n-n_0} rc - c \right) \\
&+ \sum_{n_5+n_6>0} p_4(\vec{n}) \left( \frac{n_2+n_4+n_5+1}{n-n_0} rc - c \right) \\
\pi_5 &= p_5(L)\sigma + \sum_{n_0 < n-1} p_5(\vec{n}) \left( \frac{n_2+n_4+n_5+1}{n-n_0} rc - c - \gamma_0 - \frac{\gamma_1}{n_5+n_6+1} \right) \\
\pi_6 &= p_6(L)\sigma + \sum_{n_0 < n-1} p_6(\vec{n}) \left( \frac{n_2+n_4+n_5}{n-n_0} rc - \beta - \gamma_0 - \frac{\gamma_1}{n_5+n_6+1} \right).
\end{aligned} \tag{S3}$$

In these expressions, the term

$$p_i(\vec{n}) = \frac{\binom{M_i-1}{n_i} \prod_{j \neq i} \binom{M_j}{n_j}}{\binom{M-1}{n-1}}, \tag{S4}$$

gives the probability that player  $i$  is in a group with composition  $\vec{n}$  (corresponding to the case of sampling without replacement), and  $p_i(L)$  is the shorthand notation for  $p_i(\vec{n})$  with  $\vec{n} = (n-1, 0, \dots, 0)$  (i.e., all other players abstain from the public goods game). Using the properties of multivariate hypergeometric distributions, one can simplify the payoffs in Eq. S3 (such that one does not need to sum up over all possible group compositions  $\vec{n}$ ). The simplified payoff formulas are shown in the *Appendix*.

**Derivation of the payoffs for institutions with second-order punishment.** Second-order punishment leads to a slight modification of the payoff formulas. If a punishment institution has been established, then all players who did not pay taxes (players with strategies  $S_1$  to  $S_4$ ) have to pay an additional fine  $\beta$ . Therefore, we obtain the payoffs

$$\begin{aligned}
\tilde{\pi}_i &= \pi_i & \text{for } i \in \{0, 5, 6\} \\
\tilde{\pi}_i &= \pi_i - \sum_{n_5+n_6>0} p_i(\vec{n})\beta & \text{for } i \in \{1, 2, 3, 4\},
\end{aligned} \tag{S5}$$

where the  $\pi_i$  denote the payoffs without second-order punishment given in Eq. S3. Again, these payoff formulas can be simplified (*Appendix*).

**Description of the evolutionary process.** As the basic idea of evolutionary game theory, the composition of the population is not constant; rather, individuals learn to adopt new strategies depending on the relative success of their strategies. In the following, we consider a pairwise comparison process (5). In each time step, two individuals are randomly sampled from the population: a learner and a role model. Depending on the learner's payoff  $\pi_i$  and on the role model's payoff  $\pi_j$ , the learner decides to imitate the role model's strategy with a probability that increases in the payoff difference  $\pi_j - \pi_i$ . A frequently used specification of this probability is given by the Fermi function (13)

$$f(\pi_i, \pi_j) = \frac{1}{1 + e^{-s(\pi_j - \pi_i)}} \tag{S6}$$

The parameter  $s \geq 0$  reflects the strength of selection; for  $s = 0$ , this probability is always  $1/2$ , and imitation occurs randomly. As  $s$  becomes larger, the imitation process is increasingly biased in favor of strategies that yield high payoffs. In addition to these imitation events, we allow individuals to explore the strategy space by mutating to different strategies. To implement such mutations, we assume that in each time step, a player may switch to one of the other strategies with probability  $\mu$  (with all other strategies having the same chance to be selected). In the following, we investigate the dynamics of the resulting process by performing extensive individual-based simulations. By calculating the average abundance of each strategy over time, we approximate the unique invariant distribution of the process.

**Results.** Fig. S1 shows the main results of our numerical analysis as a function of the selection strength parameter. Over the whole range of possible selection strengths, we observe that second-order punishment leads to more cooperation in the public goods game, due to a higher willingness to pay taxes. As a consequence, institutions with second-order punishment increase the average payoff of the population if the selection strength is sufficiently strong (for the parameter values of the experiment,  $s \geq 0.1$ ).

In the limit of strong selection ( $s \rightarrow \infty$ ), in which players adopt profitable strategies only, second-order punishment results in full cooperation, all players pay taxes, and average payoffs approach € 0.90 per round, as already predicted by the static equilibrium analysis of the game. On the level of individual strategies, second-order punishment leads to almost immediate fixation of righteous citizens, i.e., individuals who both pay taxes for the central authority and contribute to the common pool.

In contrast, without second-order punishment, we observe that only a quarter of individuals is willing to pay taxes, and individuals contribute to the common pool in approximately 75% of all cases, such that the resulting average payoffs are close to € 0.73, as predicted by the static model. This lower average payoff is caused by the persistence of opportunistic players, i.e., subjects who don't pay taxes and only contribute to the common pool if a punishment institution has been established (it follows that a group of opportunists yields a payoff of zero).

Thus, both modeling approaches (static and dynamic) yield the same prediction: institutions with second-order punishment are more successful in preventing tax evasion and incentivizing cooperation. Overall, second-order punishment pays in the sense that it yields higher average payoffs (despite the fact that second-order punishment institutions are slightly more costly to implement and that they decrease the payoff of tax evaders).

**Further Analysis of the Experiment. Decisions in block II.** The first two blocks of the experiment are the same for the two treatments. Thus, we aggregated the results of both treatments in Fig. 1. Indeed, for none of the 10 variables under consideration did we find a treatment effect (i.e., there were no significant differences between treatments A and B at the  $\alpha = 0.10$  level, Mann-Whitney  $U$  test with  $n_A = 10$  and  $n_B = 15$ ). Here we show the results of



each of the two treatments separately. As in the main text, we will refer to individuals that contribute to the common pool as cooperators and to individuals who do not contribute as defectors. Individuals that decide not to participate in the game will be called loners. The results were as follows:

- i) Treatment A: The addition of second-order punishment had a significantly positive effect on payoffs (payoffs increased from € 0.63 in rounds without second-order punishment to € 0.86, Wilcoxon matched-pairs signed-rank test,  $n_A = 10$ ,  $Z = 2.808$ ,  $P = 0.005$ ). Second-order punishment resulted in a significant increase of cooperation (without second-order punishment: 63.2%; with second-order punishment: 98.0%; Wilcoxon matched-pairs signed-rank test,  $Z = 2.808$ ,  $P = 0.005$ ), and significantly less defection (30.8% vs. 0.4%, Wilcoxon matched-pairs signed-rank test,  $Z = 2.808$ ,  $P = 0.005$ ). As expected, there was also a significant increase in the number of tax payers (10.8% vs. 87.8%, Wilcoxon matched-pairs signed-rank test,  $Z = 2.808$ ,  $P = 0.005$ ).
- ii) Treatment B: All results were comparable to the results of the first treatment: second-order punishment led to an increase of payoffs (from € 0.63 to € 0.86 per round, Wilcoxon matched-pairs signed-rank test,  $n_B = 15$ ,  $Z = 2.982$ ,  $P = 0.003$ ), due to significantly more cooperators (60.3% compared with 95.5%, Wilcoxon matched-pairs signed-rank test,  $Z = 3.268$ ,  $P = 0.001$ ), and less defectors (29.1% compared with 0.8%, Wilcoxon matched-pairs signed-rank test,  $Z = 3.315$ ,  $P = 0.001$ ). The fraction of tax payers increased from 8.0% to 87.6% (Wilcoxon matched-pairs signed-rank test,  $P = 0.001$ ).
- iii) Theory vs. experiment: The results in both treatments were in remarkably good agreement with the theoretical predictions. For example, the predicted average payoffs per round were € 0.73 (without second-order punishment) and € 0.90 (with second-order punishment), whereas the experiments resulted in average payoffs of € 0.63 and € 0.86, respectively. Similarly, the predicted fraction of defectors was 22.8% (without second-order punishment) and 0.0% (with second-order punishment), whereas the observed frequencies were 29.8% and 0.6%, respectively. There was a similar qualitative agreement in the number of tax payers and in the number of cooperators (Fig. S2).

**Decisions in block III.** Fig. S3 shows the individual preference for institutions with second-order punishment in the first and the last round of the third game (in case of treatment A), as well as the result of the majority vote (for treatment B). In treatment A, the players' choices were not constant over time: whereas in the first round there was a majority for institutions without second-order punishment in 8 of 10 groups, this fraction decreased to 4 of 10 groups by the end of the third block. In both cases, the hypothesis of indifference between the two punishment regimes cannot be rejected (for the first round the binomial test leads to  $P = 0.109$ , whereas for the last round we obtain  $P = 0.754$ ). In treatment B, groups showed a significant preference for a punishment institution with second-order punishment, with 12 of 15 groups having a majority for such a regime (binomial test,  $P = 0.035$ ). A Kolmogorov–Smirnov two-sample test verifies that the treatments differ significantly in the voting behavior in the beginning of block III ( $K = 1.470$ ,  $P = 0.027$ ). However, by the end of block III, the voting behavior in treatment A does not differ significantly from the outcome of the majority vote in treatment B any more (Kolmogorov–Smirnov two-sample test,  $K = 0.898$ ,  $P = 0.395$ ).

Fig. S4 shows the aggregated results of the public goods games in block III. Notably, the voting procedure and the resulting institutional choices had a systematic impact on the behavior in the actual public goods game. On average, the majority vote resulted in higher average payoffs (A: € 0.64; B: € 0.86;  $n_A = 10$ ,  $n_B = 15$ , Mann-Whitney  $U$  test,  $Z = 3.534$ ,  $P = 0.001$ ) due to a higher abundance of cooperation (A: 64.5%; B: 93.2%; Mann-

Whitney  $U$  test,  $Z = 3.528$ ,  $P = 0.001$ ) and a lower abundance of defectors (A: 22.1%; B: 5.1%; Mann-Whitney  $U$  test,  $Z = 3.421$ ,  $P = 0.001$ ). For treatment A, the overall results in block III are therefore close to the results during block II in rounds without second-order punishment. On the other hand, the results for treatment B in block III strongly reflect the characteristic features of the corresponding results during block II under an institution with second-order punishment (cf. Figs. S2 and S4).

**Instructions of the Experiment.** In the following, we provide examples of the information displayed on the subjects' laptops throughout the experiment, translated from German.

**Instructions in the beginning of the experiment. Page 1.** *Welcome to this experiment, in which you can earn money. At the beginning of this experiment, you will receive 12 Euros credited to your account. During the experiment you can win or lose money. This depends on your decisions and the decisions of the other players. Your decisions are anonymous. To ensure this, the computer assigns you a pseudonym that can be seen at the bottom left of your screen. The pseudo names are names of moons in our solar system (Ananke, Telesto, Despina, Japetus, and Kallisto). At the end of the game you will receive in cash the money in your account anonymously under your pseudo name. To render this experiment successful, it is strictly forbidden for participants to talk to each other or to communicate in any other way. After having read this text completely, please confirm by pressing the green OK-button.*

**Page 2.** *All players need to decide whether to play alone (loner) or in a group. If you decide to play alone, no further decisions have to be made. At the end of the round, the fixed amount of 0.40 Euro will be credited to your account. If you decide to play in the group, you have to make additional decisions. First, all group players are asked whether to pay taxes for policing. Thereafter, all group players are asked whether they want to invest into a group project. The decisions will be shown to you after all group members have made their decision. After that, the amount in the group project will be multiplied by 3.1 and paid to the group players in equal shares (credited to each individual account), irrespective whether they have invested into the group project or not. After having read this text completely, please confirm by pressing the green OK-button.*

**Page 3.** *First, all players are asked simultaneously: Do you want to play in the group? Answer: Yes or No. If you choose "No", you are automatically a "loner" and receive 0.40 Euro. There can be any possible mix of loners and group players. Exception: if only one player decides "group", he will become a loner automatically. After having read this text completely, please confirm by pressing the green OK button.*

**Page 4.** *Players who decided to play in the group, are asked simultaneously: Do you want to pay taxes for policing? Answer: Yes or No. (Yes pays taxes for policing, No pays nothing). The amount of the tax depends on the number of players who pay taxes:*

Number of tax payers	1	2	3	4	5
Tax per player	0.50 Euro	0.28 Euro	0.20 Euro	0.16 Euro	0.14 Euro

*After having read this text completely, please confirm by pressing the green OK button.*

At this point of the experiment, all players are informed whether someone paid taxes for the policing institution. The text that is displayed is either "Taxes have been paid" or "No taxes have been paid," followed by the sentence "After having read this text completely, please confirm by pressing the green OK button."

**Page 5.** *Thereafter, all group players are asked: Do you want to contribute to the group project? Answer: Yes or No. (Yes pays 0.50 Euros in the group project, No pays nothing). The sum of investments into the group project will be multiplied by 3.1 and will be paid to all group players in equal shares. If there was at*

least one player who previously paid taxes for policing, then all players who did not invest into the group project pay a 1.00 Euro fine. If no taxes were paid, then players who did not invest into the group project do not pay a fine. After having read this text completely, please confirm by pressing the green OK button.

**Page 6.** Examples: For a better understanding, for now the following examples only show the consequences of payments into the group project.

Pseudo name	Leda	Triton	Portia	Carpo	Galatea
Loner/group	Group	Group	Group	Group	Loner
Payment to loner	—	—	—	—	0.40
Investment into group project	0.50	0.50	0.50	0.50	—
Profit from group project	1.55	1.55	1.55	1.55	—
Total profit	1.05	1.05	1.05	1.05	0.40

Pseudo name	Leda	Triton	Portia	Carpo	Galatea
Loner/group	Group	Group	Loner	Group	Group
Payment to loner	—	—	0.40	—	—
Investment into group project	0.50	0.00	—	0.00	0.50
Profit from group project	0.78	0.78	—	0.78	0.78
Total profit	0.28	0.78	0.40	0.78	0.28

Pseudo name	Leda	Triton	Portia	Carpo	Galatea
Loner/group	Loner	Group	Group	Group	Group
Payment to loner	0.40	—	—	—	—
Investment into group project	—	0.00	0.00	0.00	0.00
Profit from group project	—	0.00	0.00	0.00	0.00
Total profit	0.40	0.00	0.00	0.00	0.00

After having read this text completely, please confirm by pressing the green OK button.

**Page 7.** Examples: The following examples additionally show the consequences of taxes for policing.

After having read this text completely, please confirm by pressing the green OK button.

Pseudo name	Leda	Triton	Portia	Carpo	Galatea
Loner/group	Group	Group	Group	Group	Loner
Payment to loner	—	—	—	—	0.40
Taxes for policing	0.50	0.00	0.00	0.00	—
Investment into group project	0.50	0.00	0.50	0.00	—
Profit from group project	0.78	0.78	0.78	0.78	—
Police penalty for noninvestment	0.00	1.00	0.00	1.00	—
Total profit	-0.22	-0.22	0.28	-0.22	0.40

Pseudo name	Leda	Triton	Portia	Carpo	Galatea
Loner/group	Group	Group	Group	Loner	Group
Payment to loner	—	—	—	0.40	—
Taxes for policing	0.00	0.28	0.00	—	0.28
Investment into group project	0.50	0.50	0.00	—	0.50
Profit from group project	1.16	1.16	1.16	—	1.16
Police penalty for noninvestment	0.00	0.00	1.00	—	0.00
Total profit	0.66	0.38	0.16	0.40	0.38

**Page 8.** The experiment starts now! You have a credit of 12 Euros on your account. After having read this text completely, please confirm by pressing the green OK button.

**Instructions after the first five rounds of the experiment. Page 1.** Police penalty for not paying taxes. From now on there is an additional consequence for group players when punishing: if you decide to pay taxes, such that your coplayers are punished for not investing then additionally those players will be punished who have invested but have not paid taxes for policing. The penalty is 1.00 euro in each case. The amount of the tax is now slightly higher and it still depends on the number of tax payers:

Number of tax payers	1	2	3	4	5
Tax per player	0.55 Euro	0.30 Euro	0.22 Euro	0.18 Euro	0.15 Euro

Example:

Pseudo name	Leda	Triton	Portia	Carpo	Galatea
Loner/group	Group	Group	Group	Loner	Group
Payment to loner	—	—	—	0.40	—
Taxes for policing	0.00	0.30	0.00	—	0.30
Investment into group project	0.50	0.50	0.00	—	0.50
Profit from group project	1.16	1.16	1.16	—	1.16
Police penalty for noninvestment	0.00	0.00	1.00	—	0.00
Police penalty for not paying taxes	1.00	0.00	1.00	—	0.00
Total profit	-0.34	0.36	-0.84	0.40	0.36

After having read this text completely, please confirm by pressing the green OK button.

**Instructions in the beginning of block II of the experiment. Page 1.** Welcome to another experiment, in which you can earn money! At the beginning of this experiment you will again receive 12 Euros credited to your account. The course of the game is exactly the same as the one you have just played. After having read this text completely, please confirm by pressing the green OK button.

**Page 2.** As in the first half of the game that you have just played: if you decide to pay taxes for policing, then all coplayers are punished for not investing into the group project. Players that do contribute to the group project but who don't pay taxes are not punished. Again, the amount of tax for policing depends on the number of tax payers:

Number of tax payers	1	2	3	4	5
Tax per player	0.50 Euro	0.28 Euro	0.20 Euro	0.16 Euro	0.14 Euro

After having read this text completely, please confirm by pressing the green OK button.

**Page 3.** Examples (which you know from the first half of the last game):

Pseudo name	Leda	Triton	Portia	Carpo	Galatea
Loner/group	Group	Group	Group	Group	Loner
Payment to loner	—	—	—	—	0.40
Taxes for policing	0.50	0.00	0.00	0.00	—
Investment into group project	0.50	0.00	0.50	0.00	—
Profit from group project	0.78	0.78	0.78	0.78	—
Police penalty for noninvestment	0.00	1.00	0.00	1.00	—
Total profit	-0.22	-0.22	0.28	-0.22	0.40

Pseudo name	Leda	Triton	Portia	Carpo	Galatea
Loner/group	Group	Group	Group	Loner	Group
Payment to loner	—	—	—	—	—
Taxes for policing	0.00	0.28	0.00	—	0.28
Investment into group project	0.50	0.50	0.00	—	0.50
Profit from group project	1.16	1.16	1.16	—	1.16
Police penalty for noninvestment	0.00	0.00	1.00	—	0.00
Total profit	0.66	0.38	0.16	0.40	0.38

After having read this text completely, please confirm by pressing the green OK button.

**Page 4.** The experiment starts now! You have a credit of 12 Euros on your account. After having read this text completely, please confirm by pressing the green OK button.

**Instructions after the first 5 rounds of block II. Page 1.** Police penalty for not paying taxes. From now on there is again an additional consequence for group players when punishing: if you decide to pay taxes, such that your coplayers are punished for not investing then additionally those players will be punished who have invested but have not paid taxes for policing. The penalty is 1.00 euro in each case. The amount of the tax is now slightly higher and it still depends on the number of tax payers:

Number of tax payers	1	2	3	4	5
Tax per player	0.55	0.30	0.22	0.18	0.15
	Euro	Euro	Euro	Euro	Euro

Example:

Pseudo name	Leda	Triton	Portia	Carpo	Galatea
Loner/group	Group	Group	Group	Loner	Group
Payment to loner	—	—	—	0.40	—
Taxes for policing	0.00	0.30	0.00	—	0.30
Investment into group project	0.50	0.50	0.00	—	0.50
Profit from group project	1.16	1.16	1.16	—	1.16
Police penalty for noninvestment	0.00	0.00	1.00	—	0.00
Police penalty for not paying taxes	1.00	0.00	1.00	—	0.00
Total profit	-0.34	0.36	-0.84	0.40	0.36

After having read this text completely, please confirm by pressing the green OK button.

**Instructions in the beginning of block III (foot-voting treatment). Page 1.** Welcome to another experiment, in which you can earn money! At the beginning of this experiment you will receive 18 Euros credited to your account. The course of the game is similar to the one you have just played. However, this time you can choose before each round whether you want to play in a community in which players are punished for not investing into the group project, or whether you want to play in a community in which additionally those coplayers are punished who invest into the public pool but who do not pay taxes. You have already played both variants of the game, but sequentially after each other and without an option to choose between them. From now on you can choose before each round whether you want to play in a community where all group members play according to one variant, or whether you want to play in a community where all group members play according to the other variant. After each round, each player sees a summary of the results of both communities. You will play many rounds. Before each round you have to decide again in

which community you want to play. After having read this text completely, please confirm by pressing the green OK button.

**Page 2.** The experiment starts now! You have a credit of 18 Euros on your account. After having read this text completely, please confirm by pressing the green OK button.

**Instructions in the beginning of block III (majority-voting treatment).** In the beginning of block III, one of the experimenters (M.M.) made the following announcement,

Before the next experiment starts, there is an election. In this election, you can decide democratically in which mode you want to play during the next experiment. For this reason, we will give you two pieces of paper, one with the text “Fine for noncontributors,” and one with the text “Fine for noncontributors and fine for tax evaders.” You already know both regimes. However, this time there will be no switch between these regimes, as in the previous experiments. Instead, the majority determines the mode for all remaining rounds of the experiment.

Then subjects had to vote by putting one of the two pieces of paper into a ballot box. Thereafter, the result of the election was announced by the experimenter by either stating “Fine for noncontributors” or “Fine for noncontributors and fine for tax evaders.” Then the corresponding computer program was started; the instructions during block III were then analogous to the previous instructions during block II.

#### Appendix: Payoffs in the Social Learning Model

Using the properties of binomial coefficients and the multinomial distribution, one can write the payoffs in Eq. S3 as follows:

$$\begin{aligned} \pi_1 &= \frac{\binom{M_0}{n-1}}{\binom{M-1}{n-1}} \sigma + B - \left( 1 - \frac{\binom{M-M_5-M_6-1}{n-1}}{\binom{M-1}{n-1}} \right) \beta \\ \pi_2 &= \frac{\binom{M_0}{n-1}}{\binom{M-1}{n-1}} \sigma + B - F_P \cdot c \\ \pi_3 &= \frac{\binom{M_0}{n-1}}{\binom{M-1}{n-1}} \sigma + B - F_N \cdot c - \left( 1 - \frac{\binom{M-M_5-M_6-1}{n-1}}{\binom{M-1}{n-1}} \right) \beta \\ \pi_4 &= \frac{\binom{M_0}{n-1}}{\binom{M-1}{n-1}} \sigma + B - (F_N + F_P) \cdot c \\ \pi_5 &= \frac{\binom{M_0}{n-1}}{\binom{M-1}{n-1}} \sigma + B_T - F_T \cdot c - \Gamma \\ \pi_6 &= \frac{\binom{M_0}{n-1}}{\binom{M-1}{n-1}} \sigma + B_T - \left( 1 - \frac{\binom{M_0}{n-1}}{\binom{M-1}{n-1}} \right) \beta - \Gamma, \end{aligned} \quad [S7]$$

with

$$\begin{aligned}
 B &= rc \left( 1 - \frac{M}{n(M-M_0)} + \frac{\binom{M_0}{n}}{\binom{M-1}{n-1}} \frac{1}{M-M_0} \right) \frac{(M_4+M_2)(M-M_0-M_6-1) + (M_1+M_3-1)M_5}{(M-M_0-M_6-M_5-1)(M-M_0-1)} \\
 &\quad + rc \left( 1 - \frac{M-M_5-M_6}{n(M-M_5-M_6-M_0)} + \frac{\binom{M_0}{n}}{\binom{M-1}{n-1}} \frac{1}{M-M_5-M_6-M_0} \right) \times \frac{\binom{M-M_5-M_6-1}{n-1}}{\binom{M-1}{n-1}} \frac{M_3-M_2}{M-M_0-M_5-M_6-1} \\
 F_P &= 1 - \frac{\binom{M-M_5-M_6-1}{n-1}}{\binom{M-1}{n-1}} - \frac{r}{M-M_0-1} \left( \frac{M}{n} - 1 - \frac{\binom{M_0}{n}}{\binom{M-1}{n-1}} \right) \\
 &\quad + \frac{r}{M-M_0-M_5-M_6-1} \frac{\binom{M-M_5-M_6-1}{n-1}}{\binom{M-1}{n-1}} \times \left( \frac{M-M_5-M_6}{n} - 1 - \frac{\binom{M_0}{n}}{\binom{M-M_5-M_6-1}{n-1}} \right) \\
 F_N &= \frac{\binom{M-M_5-M_6-1}{n-1} - \frac{\binom{M_0}{n-1}}{\binom{M-1}{n-1}}}{\binom{M-1}{n-1}} + \frac{\binom{M_0}{n-1}}{\binom{M-1}{n-1}} r \\
 &\quad - \frac{r}{M-M_0-M_5-M_6-1} \frac{\binom{M-M_5-M_6-1}{n-1}}{\binom{M-1}{n-1}} \times \left( \frac{M-M_5-M_6}{n} - 1 - \frac{\binom{M_0}{n}}{\binom{M-M_5-M_6-1}{n-1}} \right) \\
 B_T &= rc \left( 1 - \frac{M}{n(M-M_0)} + \frac{\binom{M_0}{n}}{\binom{M-1}{n-1}} \frac{1}{M-M_0} \right) \frac{M_4+M_5+M_2}{M-M_0-1} \\
 F_T &= 1 - \frac{r}{M-M_0-1} \left( \frac{M}{n} - 1 \right) + \frac{\binom{M_0}{n-1}}{\binom{M-1}{n-1}} \left( \frac{r}{n} \frac{M_0+1}{M-M_0-1} + r \frac{M-M_0-2}{M-M_0-1} - 1 \right) \\
 \Gamma &= \left( 1 - \frac{\binom{M_0}{n-1}}{\binom{M-1}{n-1}} \right) \gamma_0 + \frac{\gamma_1 M}{n(M_5+M_6)} \left( 1 - \frac{\binom{M-M_5-M_6}{n} + (M_5+M_6) \binom{M_0}{n-1}}{\binom{M}{n}} \right).
 \end{aligned}$$

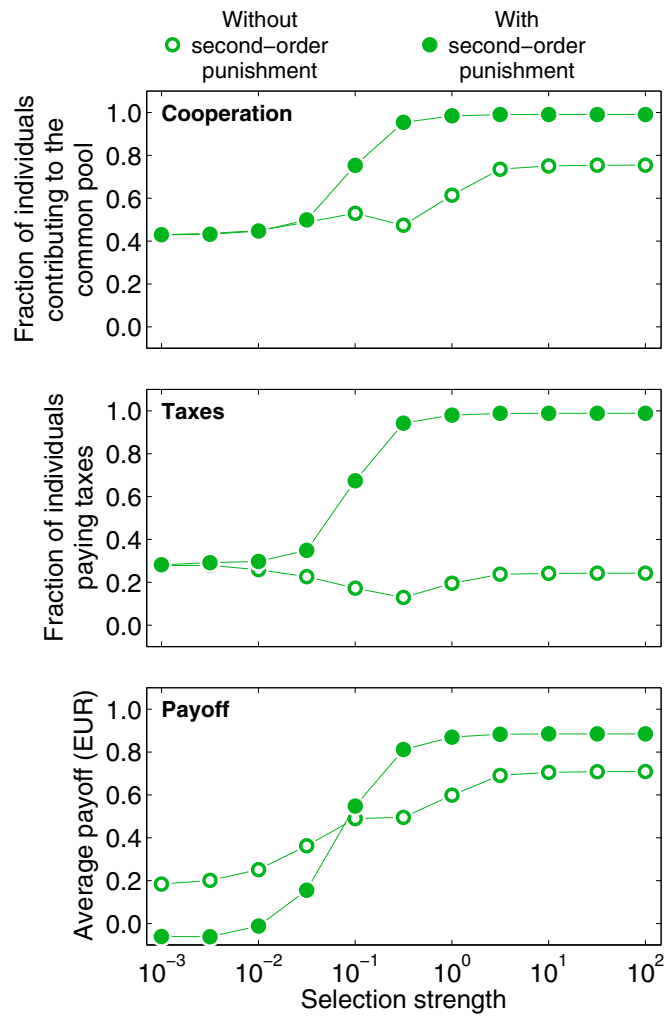
[S8]

The variable  $B$  corresponds to the expected benefit that a non-contributing player derives from the common pool.  $F_P$  denotes the effective marginal cost of contributing when a punishment institution has been established, whereas  $F_N$  gives the marginal cost when no such institution has been established.  $B_T$  is the expected benefit that a tax-payer derives from the common pool, with  $F_T$  being the effective cost of contributing to the common pool when being a tax payer. Last,  $\Gamma$  gives the expected amount that a tax payer needs to pay to the central authority.

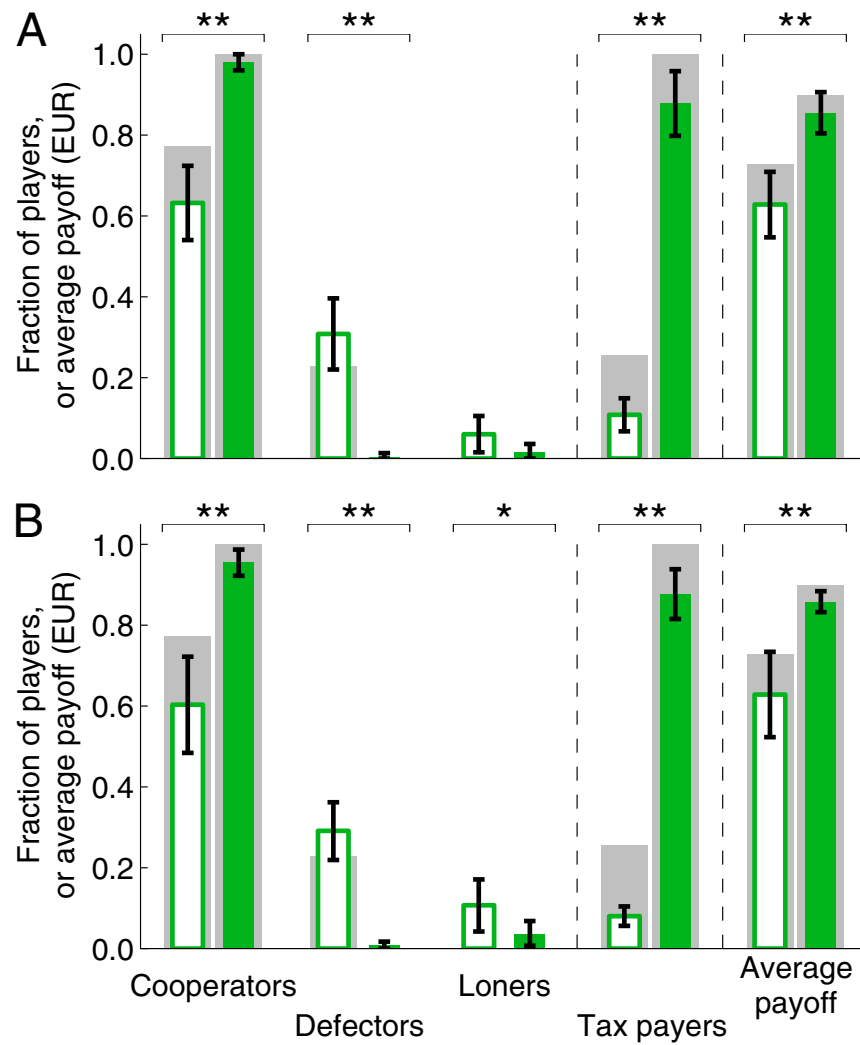
To calculate the expected payoffs in the case with second-order punishment, according to Eq. S5, we note that for  $i \in \{1, 2, 3, 4\}$

$$\sum_{n_5+n_6>0} p_i(\vec{n})\beta = \left( 1 - \frac{\binom{M-M_5-M_6-1}{n-1}}{\binom{M-1}{n-1}} \right) \beta. \quad [\text{S9}]$$

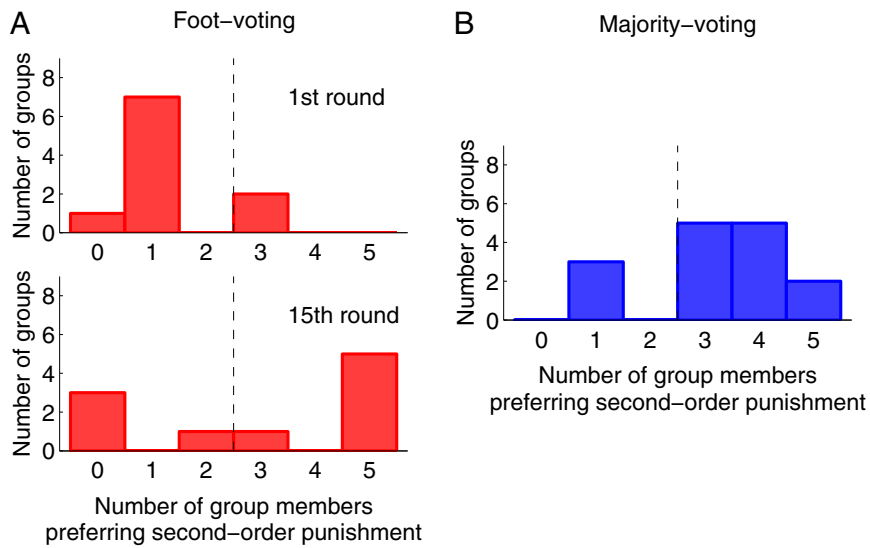
1. Semmann D, Krambeck HJ, Milinski M (2005) Reputation is valuable within and outside one's own social group. *Behav Ecol Sociobiol* 57(6):611–616.
2. Selten R (1965) Spieltheoretische Behandlung eines Oligopolmodells mit Nachfrageträgheit. *Zeitschrift Gesamte Staatswissenschaft* 121(2):301–324.
3. Fudenberg D, Tirole J (1998) *Game Theory* (MIT Press, Cambridge, MA), 6th Ed.
4. Cressman R, Schlag KH (1998) The dynamic (in)stability of backwards induction. *J Econ Theory* 83(2):260–285.
5. Traulsen A, Nowak MA, Pacheco JM (2006) Stochastic dynamics of invasion and fixation. *Phys Rev E Stat Nonlin Soft Matter Phys* 74(1 Pt 1):011909.
6. Sigmund K, De Silva H, Traulsen A, Hauert C (2010) Social learning promotes institutions for governing the commons. *Nature* 466(7308):861–863.
7. Diekmann A (1985) Volunteer's dilemma. *J Conflict Resolut* 29(4):605–610.
8. Nowak MA (2006) *Evolutionary Dynamics* (Harvard Univ Press, Cambridge, MA).
9. Sigmund K (2010) *The Calculus of Selfishness* (Princeton Univ. Press, Princeton, NJ).
10. Traulsen A, Hauert C (2009) Stochastic evolutionary game dynamics. *Reviews of Nonlinear Dynamics and Complexity*, ed Schuster HG (Wiley-VCH, Weinheim), Vol II, pp 25–61.
11. Hauert C, Traulsen A, Brandt H, Nowak MA, Sigmund K (2007) Via freedom to coercion: the emergence of costly punishment. *Science* 316(5833):1905–1907.
12. Schoenmakers S (2013) Pool-punishment and opportunistic cooperation in voluntary and compulsory games. An evolutionary game theory model. MS thesis (Univ of Oldenburg, Oldenburg, Germany).
13. Blume LE (1993) The statistical mechanics of strategic interaction. *Games Econ Behav* 5(3):387–424.



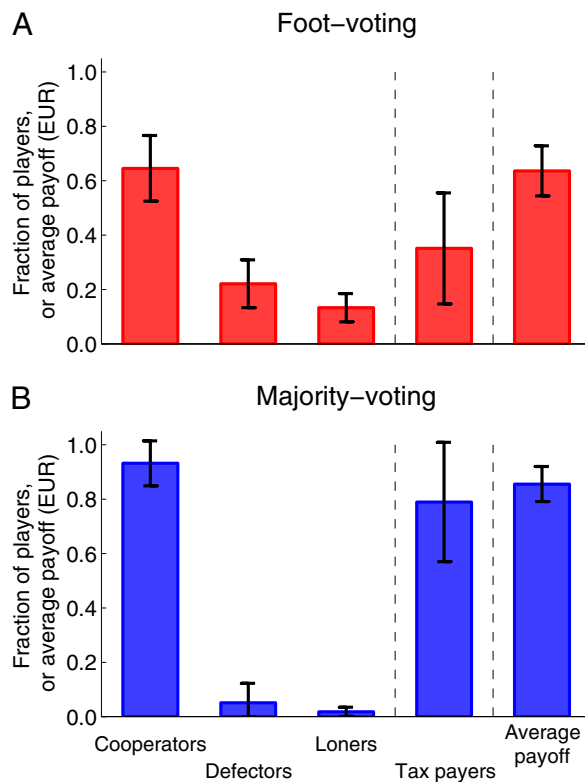
**Fig. S1.** Individual-based simulations show that institutions with second-order punishment lead to more cooperation, less tax evasion, and higher payoffs. For each time step of the simulation, we recorded the fraction of tax payers (subjects applying strategy  $S_5$  or  $S_6$ ), the fraction of contributors (either unconditionally by playing  $S_4$  or  $S_5$ , or conditionally, when applying strategy  $S_2$  or  $S_3$ , weighted by the probability that these subjects contribute), and the resulting average payoff. Simulations are run for  $10^7$  time steps, with the parameters used in the experiment, population size  $M=100$ , and mutation rate  $\mu=0.01$ .



**Fig. 52.** An analysis of the results in block II, separately for each treatment. The bars show the average payoff per round (in Euros) and the fraction of players choosing a given behavior. Empty bars correspond to the outcome in rounds without second-order punishment, and filled bars show the rounds with second-order punishment. \*Significant at the  $\alpha = 0.05$  level; \*\*significant at the  $\alpha = 0.01$  level (using Wilcoxon matched-pairs signed-rank tests). Gray bars represent the predicted value based on the theoretical model. (A) Results for the foot voting treatment, and (B) results of the majority voting treatment. As expected, there were no significant differences between these two treatments during block II.



**Fig. 53.** Voting behavior in block III across the different treatments. Each graph shows the number of groups where a given number of subjects voted for second-order punishment. (A) For the foot voting treatment we show the players’ decisions in the first and the last round. (B) For the majority voting treatment, the graph shows the players’ decisions in the election before block III.



**Fig. 54.** Aggregated results of the public goods games in block III. As in Fig. S1, bars show the average payoff per round (in Euros) and the fraction of players choosing a given behavior. Averages are taken over all rounds of the third block, with error bars depicting the 95% CI. (A) Results for the foot voting treatment are qualitatively similar to the results of block II during periods without second-order punishment. (B) In contrast, outcomes in the majority voting treatment resemble the respective findings of block II during periods with second-order punishment.



**Table S1. Parameters of the public goods game**

Parameter	Description	Parameter value in the experiment
$n$	Group size for the public goods game	$n = 5$
$\sigma$	Secure payoff when abstaining from the public goods game	$\sigma = \text{€ } 0.40$
$c$	Cost of contributing to the common pool	$c = \text{€ } 0.50$
$r$	Multiplication factor for contributions to the common pool	$r = 3.1$
$\beta$	Punishment fine for noncontributors (and tax evaders) in case that a punishment institution is established	$\beta = \text{€ } 1.00$
$\gamma_0 + \gamma_1/i$	Taxes for the punishment institution, as a function of the number of tax payers $i$	$\gamma_0 = \text{€ } 0.05$ ; $\gamma_1 = \text{€ } 0.45$ (without 2OP) $\gamma_1 = \text{€ } 0.50$ (with 2OP)

We assume that the public goods game is a social dilemma,  $1 < r < n$ , and punishment acts as a deterrent,  $\beta > c$  and  $\beta > \gamma_0 + \gamma_1/i$  for all  $i \geq 1$ .

**Table S2. Possible strategies in the public goods game**

Strategy	Description
$S_0$	Decides to abstain from the public goods game
$S_1$	Participates in the game, but does not pay taxes and does not contribute to the common pool
$S_2$	Participates in the game, but does not pay taxes and only contributes to the common pool if someone paid taxes
$S_3$	Participates in the game, but does not pay taxes and only contributes to the common pool if no one paid taxes
$S_4$	Participates in the game, does not pay taxes, but always contributes to the common pool
$S_5$	Participates in the game, pays taxes, and contributes to the common pool
$S_6$	Participates in the game, pays taxes, but does not contribute to the common pool