

Social learning promotes institutions for governing the commons

Karl Sigmund^{1,2}, Hannelore De Silva³, Arne Traulsen⁴ & Christoph Hauert⁵

Theoretical and empirical research highlights the role of punishment in promoting collaborative efforts^{1–5}. However, both the emergence and the stability of costly punishment are problematic issues. It is not clear how punishers can invade a society of defectors by social learning or natural selection, or how second-order free-riders (who contribute to the joint effort but not to the sanctions) can be prevented from drifting into a coercion-based regime and subverting cooperation. Here we compare the prevailing model of peer-punishment^{6–8} with pool-punishment, which consists in committing resources, before the collaborative effort, to prepare sanctions against free-riders. Pool-punishment facilitates the sanctioning of second-order free-riders, because these are exposed even if everyone contributes to the common good. In the absence of such second-order punishment, peer-punishers do better than pool-punishers; but with second-order punishment, the situation is reversed. Efficiency is traded for stability. Neither other-regarding tendencies or preferences for reciprocity and equity, nor group selection or prescriptions from higher authorities, are necessary for the emergence and stability of rudimentary forms of sanctioning institutions regulating common pool resources and enforcing collaborative efforts.

Many economic experiments on ‘public goods games’ (PGGs) have shown that a substantial fraction of players are willing to incur costs to impose fines on exploiters, that is, those who do not contribute to the joint effort^{1–8}. As a consequence, the threat of punishment looms credibly enough to increase the average level of pro-social contributions. However, the sanctioning system is itself a public good. Thus, punishers are often seen as altruistic, because others benefit from their costly efforts^{9–13}. Conversely, those who refrain from punishing exploiters are ‘second-order free-riders’. Among self-interested agents, second-order free-riding should spread and ultimately cause the collapse of cooperation.

A solution is to punish second-order free-riders also¹⁴. But such ‘second-order punishment’ risks being subverted by third-order free-riders in turn, leading to infinite regress. Moreover, if everyone contributes to the public good, second-order free-riders will not be spotted. Their number can grow through neutral drift, ultimately allowing defectors to invade with impunity. We show how a simple mechanism can overcome this problem.

There exist a variety of sanctioning systems. Most experiments on public goods with punishment have considered peer-punishment: after the PGG, individuals can impose fines on exploiters, at a cost to themselves. Interestingly, the first experiment on public goods with punishment¹⁵ considered a different mechanism: players decide whether to contribute to a ‘punishment pool’ before contributing to the public goods. This can be viewed as a first step towards an institutionalized mechanism for punishing exploiters, and compared with the self-financed contract enforcement games in *Governing the*

*Commons*¹⁶. It is like paying towards a police force, whereas peer-punishers take law enforcement into their own hands.

Peer- and pool-punishment are both expensive ways to impose negative incentives on free-riders. In many economic experiments, the increase in cooperation is more than matched by the costs of punishment, and an overall reduction of total pay-off is observed^{18,9}. Because the costs of pool-punishment arise even when there are no exploiters to be punished, it seems even more socially expensive than peer-punishment. However, the issue of second-order punishment favours pool-punishment. If everyone contributes to the public good, then peer-punishers are not distinguishable from second-order free-riders. By contrast, pool-punishers declare themselves beforehand. We may expect that pool-punishment leads more easily to a second-order punishment regime and, hence, to more stability.

Because sanctioning institutions, as known from social history, usually forbid individuals to take the law into their own hands, it is also worthwhile to investigate the competition between peer- and pool-punishment. A model based on evolutionary game theory shows that both peer- and pool-punishment can emerge, if participation in the joint effort is optional rather than compulsory. Pool-punishment requires second-order punishment, whereas peer-punishment is little affected by it. Both sanctioning mechanisms can evolve if players simply imitate whatever yields the highest pay-off. If peer-punishers compete with pool-punishers, all depends on second-order punishment. Without it, the population is dominated by peer-punishers. With it, pool-punishers take over, although the average income is thereby reduced.

A ‘punishment fund’ can be viewed as a rudimentary institution to uphold the common interest. Many small-scale societies use this principle, for instance by hiring an enforcer. In *Governing the Commons*¹⁶, several examples of self-financed contract enforcement are described. They concern the provisioning and the appropriation of common resources, for instance high mountain meadows (the ‘commons’), irrigation systems or inshore fisheries. Our model shows that individuals can spontaneously adopt a self-governing institution to monitor contributions and sanction free-riders. It needs no top-down prescriptions from higher authorities, nor great feats of planning: trial and error, and the imitation of successful examples, can lead to a social contract among individuals guided by self-interest.

To model a PGG, we assume that if $N \geq 2$ individuals participate in the interaction, each can decide whether to contribute a fixed amount, $c > 0$, to the common pool. This amount will be multiplied by a factor of $r > 1$ and then divided among the $N - 1$ other players. If all contribute, they obtain $(r - 1)c$ each. Because contributors do not benefit from their own contribution, self-interested players ought to contribute nothing. If all do this, their pay-off will be zero. This reveals a social dilemma.

¹Faculty of Mathematics, University of Vienna, A-1090 Vienna, Austria. ²International Institute for Applied Systems Analysis, A-2361 Laxenburg, Austria. ³WU (Vienna University of Economics and Business), A-1090 Vienna, Austria. ⁴Max Planck Institute for Evolutionary Biology, 24306 Ploen, Germany. ⁵Department of Mathematics, University of British Columbia, Vancouver, British Columbia V6T 1Z2, Canada.

Pool-punishers not only contribute c to the PGG, but also, beforehand, an amount, G , to a punishment pool. Free-riders will be fined an amount, BN_v , proportional to the number, N_v , of pool-punishers. In the case of second-order punishment, second-order free-riders will be fined the same amount. Peer-punishers contribute c to the PGG, and after the game impose a fine, β , on each free-rider in their group, at a cost γ . If N_w peer-punishers are in the group, each defector pays a total fine βN_w . In case of second-order punishment, second-order defectors are treated just like defectors.

Let us assume that the game is not compulsory^{11,17}. Some players may abstain from the joint enterprise. They can do something else instead, and earn a pay-off, σ , independent of what others are doing. If only one player is willing to engage in the joint effort, there will be no PGG and the solitary would-be participant also earns σ .

Let M denote the population size; X the number of players who participate in the PGG and contribute, but do not punish; Y the number of defectors, who participate but contribute neither to the PGG nor to the sanctions; Z the number of non-participants; V the number of pool-punishers; and W the number of peer-punishers. Random samples of N individuals are faced with the opportunity of a joint enterprise. Social learning leads to preferential copying of successful strategies. We obtain their long-run frequencies by numerical simulations (compare with Figs 1, 2 and 3). In a limiting case, we obtain analytic results (Supplementary Information) that we now describe.

Let us first neglect peer-punishment, and assume that the pay-off, σ , for non-participants lies between zero (obtained if all free-ride) and $(r-1)c-G$ (obtained if all contribute to the PGG and the punishment pool). The inequality

$$0 < \sigma < (r-1)c - G \quad (1)$$

highlights that participating in the joint enterprise is a venture that succeeds if most participants contribute and fails if most do not.

In the absence of second-order punishment, the long-run frequencies in the (X, Y, Z, V) subpopulations are $(2, 2, 2, 1)/7$ and little cooperation is achieved. With second-order punishment, the corresponding long-run frequencies are $(0, 0, 0, 1)$. The population is dominated by pool-punishers enforcing cooperation. If the game is compulsory (that is, $Z = 0$), the population consists of free-riders only.

Alternatively, if we neglect pool-punishment, and assume that

$$0 < \sigma < (r-1)c \quad (2)$$

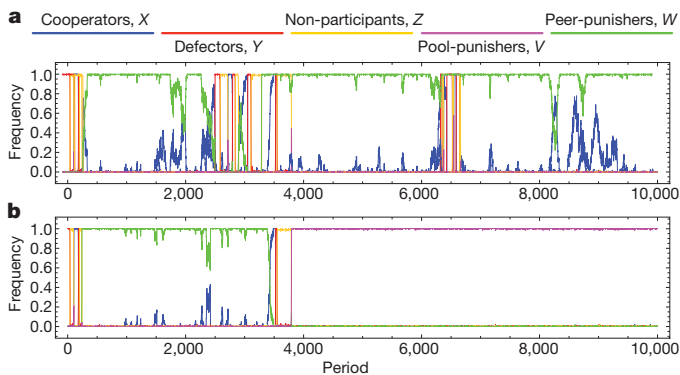


Figure 1 | Time evolution of the competition between peer-punishment and pool-punishment. Two typical individual-based simulation runs, without (a) and with (b) second-order punishment. In a, peer-punishers prevail most of the time, but sometimes second-order free-riders invade. In this case, defectors and then non-participants take over before peer-punishment is re-established. In b, pool-punishers eventually establish a very stable regime. Parameters: $N = 5$, $r = 3$, $c = 1$, $\sigma = 1$, $\gamma = \beta = 0.7$, $B = G = 0.7$, $M = 100$ and $\mu = 10^{-3}$. The updating is by strong imitation ($s \rightarrow +\infty$); that is, players with lower average pay-off always imitate players with higher average pay-off. The initial population consists of defectors only.

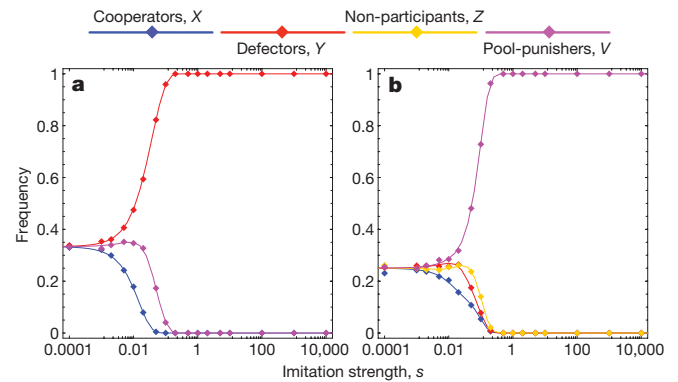


Figure 2 | Pool-punishment in compulsory and voluntary PGGs. Results of extensive simulations based on social learning (Supplementary Information). To obtain reliable average frequencies, each player updates 10^7 times. Data points are supported by analytical approximations (solid lines). Parameters are as in Fig. 1, but with $\mu = 10^{-6}$ and with variable imitation strength, s . For small s values, updating is mostly random and frequencies of all strategies are roughly equal. Discrimination between strategies increases with s . a, Compulsory PGGs lead for larger s values to a regime of defectors. b, In voluntary PGGs, the cycle $X \rightarrow Y \rightarrow Z \rightarrow X$ provides an escape from the defectors' regime through recurrent opportunities to establish a sanctioning system with second-order punishment.

the long-run frequencies in the (X, Y, Z, W) subpopulations are $(2, 2, 2, M+2)/(M+8)$ and punishers prevail, with or without second-order punishment. Again, if the game is compulsory, only free-riders survive in the long run.

In the competition between peer- and pool-punishers without second-order punishment, peer-punishers win. The long-run frequencies in the (X, Y, Z, V, W) subpopulations are $(6, 6, 4, 1, 3M+6)/(3M+23)$. With second-order punishment, pool-punishers win, and the corresponding frequencies are $(0, 0, 0, 1, 0)$.

Repression of free-riding is a basic theme for several major transitions in evolution¹⁸, and can lead to evolutionarily stable strategies allocating part of the contribution towards suppressing competition¹⁹. In human societies, sanctions are ubiquitous^{4,16,20,21}. Peer-punishment emerges more easily than pool-punishment, because it requires no second-order punishment, and inequality (2) is weaker than inequality (1). But with second-order sanctions, pool-punishment out-competes peer-punishment, despite being socially expensive. Both types of punishment only emerge, in our model, if players can opt out of the joint enterprise. This restricts the range of applications^{22,23}. However, there is considerable evidence that cooperation

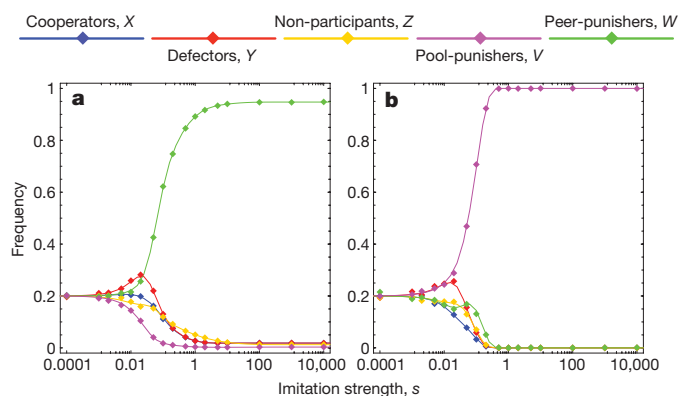


Figure 3 | The competition between peer- and pool-punishers in voluntary PGGs. a, Without second-order punishment, peer-punishers prevail but cooperation keeps breaking down and needs to be re-established (compare with Fig. 1a). b, With second-order punishment, pool-punishers prevail. Simulations and parameters are as in Fig. 2.

can increase, if participation is voluntary rather than compulsory^{24–26} (see Supplementary Information for an intuitive explanation).

Many early experiments on public goods with punishment terminated after six or ten rounds, and although punishment usually increased the propensity to cooperate, the overall income was often less than without punishment^{2,8,9}. But if the number of rounds is sufficiently large, cooperation becomes common³. As long as players avoid antisocial punishment of contributors⁵ (a feature not included in our model), peer-punishment becomes cost free. Pool-punishment entails fixed costs and thus is less efficient. However, peer-punishment is ill-suited for second-order punishment, as has also been observed empirically²⁷. Pool-punishment is more conducive to second-order punishment. A sanctioning institution should view anyone not contributing to its upkeep as a defector and resort to second-order punishment. Adding second-order punishment may add to the cost of sanctioning, but as long as inequality (1) holds, the results are unaffected.

Experimental PGGs allowing players to opt, from round to round, between treatments with or without peer-punishment²⁸, or to vote on whether to forbid antisocial punishment²⁹, suggest intermediary stages towards pool-punishment. Further steps towards endogenous institution formation are analysed in refs 23, 30. We considered players motivated entirely by self-interest, and did not assume preferences for reciprocity or equity²¹. This obviously does not mean that such preferences do not exist. Their emergence may actually have been favoured by the prevalence of sanctioning institutions over thousands of years.

We left out many important issues, such as quorum-sensing and signalling, reputation and opportunism, repeated interactions and graduated punishment, and did not specify how pool-punishment is actually set up. Our model is minimalistic, but allows proof of principle. Origins of institutions are notoriously difficult to trace, but we have shown that they can emerge spontaneously among self-interested individuals.

METHODS SUMMARY

We apply evolutionary game theory to populations of fixed size, M , and variable composition, X , Y , Z , V and W (the numbers of players using the five strategies for the optional PGG with peer- or pool-punishment). We compute the pay-offs obtained by players using these strategies. The pay-off differences define the probabilities that the strategies are copied through social learning, as a function of a parameter, $s \geq 0$, measuring 'imitation strength'. Together with an 'exploration rate', $\mu \geq 0$, which specifies the propensity to switch randomly to another strategy, this defines a stochastic process describing the evolution of the frequencies X , Y , Z , V and W . We compute their stationary distributions (which correspond to the relative frequencies in the long run) both numerically and, in a limiting case, analytically, and check these values by individual-based simulations. This allows us to compare the evolution of any subset of the five strategies under social learning. For further details, see Supplementary Information.

Received 12 March; accepted 24 May 2010.

Published online 14 July 2010.

1. Fehr, E. & Gächter, S. Cooperation and punishment in public good experiments. *Am. Econ. Rev.* **90**, 980–994 (2000).
2. Rockenbach, B. & Milinski, M. The efficient interaction of indirect reciprocity and costly punishment. *Nature* **444**, 718–723 (2006).
3. Gächter, S., Renner, E. & Sefton, M. The long-run benefits of punishment. *Science* **322**, 1510–1512 (2008).
4. Henrich, J. *et al.* Costly punishment across human societies. *Science* **312**, 1767–1770 (2006).

5. Herrmann, B., Thoenig, C. & Gächter, S. Antisocial punishment across societies. *Science* **319**, 1362–1367 (2008).
6. Fehr, E. & Gächter, S. Altruistic punishment in humans. *Nature* **415**, 137–140 (2002).
7. Gardner, A. & West, S. A. Cooperation and punishment, especially in humans. *Am. Nat.* **164**, 753–764 (2004).
8. Egas, M. & Riedl, A. The economics of altruistic punishment and the maintenance of cooperation. *Proc. R. Soc. B* **275**, 871–878 (2008).
9. Fehr, E. & Rockenbach, B. Detrimental effects of sanctions on human altruism. *Nature* **422**, 137–140 (2003).
10. Boyd, R., Gintis, H., Bowles, S. & Richerson, P. The evolution of altruistic punishment. *Proc. Natl Acad. Sci. USA* **100**, 3531–3535 (2003).
11. Fowler, J. H. Altruistic punishment and the origin of cooperation. *Proc. Natl Acad. Sci. USA* **102**, 7047–7049 (2005).
12. Nakamaru, M. & Iwasa, Y. The evolution of altruism and punishment: role of the selfish punisher. *J. Theor. Biol.* **240**, 475–488 (2006).
13. Lehmann, L., Rousset, F., Roze, D. & Keller, L. Strong reciprocity or strong ferocity? A population genetic view of the evolution of altruistic punishment. *Am. Nat.* **170**, 21–36 (2007).
14. Boyd, R. & Richerson, P. J. Punishment allows the evolution of cooperation (or anything else) in sizable groups. *Ethol. Sociobiol.* **13**, 171–195 (1992).
15. Yamagishi, T. The provision of a sanctioning system as a public good. *J. Pers. Soc. Psychol.* **51**, 110–116 (1986).
16. Ostrom, E. *Governing the Commons: The Evolution of Institutions for Collective Action* (Cambridge Univ. Press, 1990).
17. Hauert, C., Traulsen, A., Nowak, M. A., Brandt, H. H., & Sigmund, K. Via freedom to coercion: the emergence of costly punishment. *Science* **316**, 1905–1907 (2007).
18. Maynard Smith, J. & Szathmari, E. *The Major Transitions in Evolution* (Oxford Univ. Press, 1997).
19. Frank, S. A. Mutual policing and repression of competition in the evolution of cooperative groups. *Nature* **377**, 520–522 (1995).
20. Levin, S. A. (ed.) *Games, Groups, and the Global Good* (Springer, 2009).
21. Falk, A., Fehr, E. & Fischbacher, U. in *The Drama of the Commons* (eds Ostrom, L. *et al.*) 157–191 (National Academy, 2002).
22. Mathew, S. & Boyd, R. When does optional participation allow the evolution of cooperation? *Proc. R. Soc. B* **276**, 1167–1174 (2009).
23. Boyd, R., Gintis, H. & Bowles, S. Coordinated punishment of defectors sustains cooperation and can proliferate when rare. *Science* **328**, 617–620 (2010).
24. Orbell, J. H. & Dawes, R. M. Social welfare, cooperator's advantage, and the option of not playing the game. *Am. Sociol. Rev.* **58**, 787–800 (1993).
25. Hauert, C., De Monte, S., Hofbauer, J. & Sigmund, K. Volunteering as a Red Queen mechanism for cooperation. *Science* **296**, 1129–1132 (2002).
26. Semmann, D., Krambeck, H. J. & Milinski, M. Volunteering leads to rock-paper-scissors dynamics in a public goods game. *Nature* **425**, 390–393 (2003).
27. Kiyonari, T., Barclay, P., Wilson, M. & Daly, D. Second order punishment in one-shot prisoner's dilemma. *Int. J. Psychol.* **39**, 329–334 (2004).
28. Güreker, O., Irlenbush, B. & Rockenbach, B. The competitive advantage of sanctioning institutions. *Science* **312**, 108–111 (2006).
29. Ertan, A., Page, T. & Putterman, L. Who to punish? Individual decisions and majority rule in mitigating the free rider problem. *Eur. Econ. Rev.* **53**, 495–511 (2009).
30. Kosfeld, M., Riedl, A. & Okada, A. Institution formation in public goods games. *Am. Econ. Rev.* **99**, 1335–1355 (2009).

Supplementary Information is linked to the online version of the paper at www.nature.com/nature.

Acknowledgements K.S. acknowledges TECT I-104 G15, A.T. thanks the Emmy Noether programme of the DFG and C.H. thanks NSERC (Canada).

Author Contributions All authors were involved in the design and analysis of the model. H.D.S. and C.H. ran the simulations, A.T. and C.H. did the numerical analysis, and K.S. wrote the paper.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of this article at www.nature.com/nature. Correspondence and requests for materials should be addressed to K.S. (karl.sigmond@univie.ac.at).