# Runaway selection for cooperation and strict-and-severe punishment

Mayuko Nakamaru [a,*], Ulf Dieckmann [b]

[a] Department of Value and Decision Science, Tokyo Institute of Technology, 2-12-1-W9-35, O-okayama, Meguro-ku, Tokyo 152-8552, Japan
[b] Evolution and Ecology Program, International Institute for Applied Systems Analysis, Schlossplatz 1, A-2361 Laxenburg, Austria

## ABSTRACT

Punishing defectors is an important means of stabilizing cooperation. When levels of cooperation and punishment are continuous, individuals must employ suitable social standards for defining defectors and for determining punishment levels. Here we investigate the evolution of a social reaction norm, or psychological response function, for determining the punishment level meted out by individuals in dependence on the cooperation level exhibited by their neighbors in a lattice-structured population. We find that (1) cooperation and punishment can undergo runaway selection, with evolution towards enhanced cooperation and an ever more demanding punishment reaction norm mutually reinforcing each other; (2) this mechanism works best when punishment is strict, so that ambiguities in defining defectors are small; (3) when the strictness of punishment can adapt jointly with the threshold and severity of punishment, evolution favors the strict-and-severe punishment of individuals who offer slightly less than average cooperation levels; (4) strict-and-severe punishment naturally evolves and leads to much enhanced cooperation when cooperation without punishment would be weak and neither cooperation nor punishment are too costly; and (5) such evolutionary dynamics enable the bootstrapping of cooperation and punishment, through which defectors who never punish gradually and steadily evolve into cooperators who punish those they define as defectors.

© 2008 Elsevier Ltd. All rights reserved.

## 1. Introduction

Understanding the evolution of cooperation is one of the greatest challenges in evolutionary biology and the social sciences. Even though several general mechanisms are widely recognized to facilitate the emergence and maintenance of cooperation (as reviewed, e.g., by Nowak, 2006), many questions of a more detailed nature are still unresolved. Kin selection (Hamilton, 1964) explains the evolution of altruism among relatives. Direct reciprocity in repeated interactions (Axelrod and Hamilton, 1981) and indirect reciprocity enabled by reputation dynamics (e.g., Nowak and Sigmund, 1998; Leimar and Hammerstein, 2001; Panchanathan and Boyd, 2003; Brandt and Sigmund, 2004; Ohtsuki and Iwasa, 2004; Nakamaru and Kawata, 2004; Takahashi and Mashima, 2006) promote the evolution of cooperation among non-relatives. Group selection (e.g., Sober and Wilson, 1998) and selection shaped by local interactions (e.g., Matsuda, 1987; Nowak and May, 1992; Nakamaru et al., 1997, 1998; Le Galliard et al., 2003, 2005; Ohtsuki et al., 2006) may

advance cooperation in ways that can often be interpreted as generalizations of kin selection (Lehmann et al., 2007a).

Cooperation is promoted by the punishment of defectors (Axelrod, 1986; Boyd and Richerson, 1992; Clutton-Brock and Parker, 1995; Henrich and Boyd, 2001; Rockenbach and Milinski, 2006; Sigmund, 2007), and so-called altruistic punishment occurs when the direct costs of punishing are outweighed by the indirect benefits of such behavior (Yamagishi, 1986; Gintis, 2000; Sigmund et al., 2001; Fehr and Gächter, 2002; Boyd et al., 2003; Fehr and Rockenbach, 2003; Bowles and Gintis, 2004; Fehr and Fischbacher, 2004a; Gardner and West, 2004; Shinada et al., 2004; Fowler, 2005; Nakamaru and Iwasa, 2005, 2006; Brandt et al., 2006; Henrich et al., 2006; Eldakar et al., 2007; Hauert et al., 2007; Lehmann et al., 2007b; Eldakar and Wilson, 2008).

In this study, we investigate the evolution of a social reaction norm, or psychological response function, for punishment. This norm determines the threshold of encountered cooperation below which individuals punish, how strictly they apply such a threshold, and how severely they punish when they do so. In addition, we allow individuals to choose their level of cooperation from a continuum of strategies (Doebeli and Knowlton, 1998; Roberts and Sherratt, 1998; Wahl and Nowak, 1999a, b; Killingback et al., 1999; Killingback and Doebeli, 2002; Le Galliard et al., 2003, 2005; Doebeli et al., 2004). In this way, we examine the joint evolution of four continuous strategies determining, respectively,

* Corresponding author. Tel.: +81 3 5734 3365; fax: +81 3 5734 3618.
 E-mail addresses: nakamaru@valdes.titech.ac.jp (M. Nakamaru), dieckmann@iiasa.ac.at (U. Dieckmann).

the cooperation level and the threshold, strictness, and severity of punishment. Among other questions, this allows us to appraise the potential for selfish punishment and strong reciprocity: selfish punishers do not cooperate but nevertheless punish non-cooperators, whereas strong reciprocators cooperate and punish non-cooperators. Our analysis of joint evolution also allows us to compare our results with a preceding theoretical study suggesting that in a metapopulation setting the joint evolution of cooperation and punishment leads to the collapse of cooperation unless cooperation and punishment are perfectly linked traits (Lehmann et al., 2007b).

Viscous populations, exhibiting local interactions on a lattice or a more general social network, have been shown to promote the evolution of continuous cooperation strategies (Killingback et al., 1999; Le Galliard et al., 2003, 2005), as well as the joint evolution of discrete strategies of cooperation and punishment (Brandt et al., 2003; Nakamaru and Iwasa, 2005, 2006). Our study extends this earlier work to the joint and gradual evolution of continuous strategies of cooperation and punishment. In this wider context, we examine adaptable social reaction norms for punishment, analyzing their evolutionary determinants and consequences.

## 2. Methods

We consider populations in which individuals occupy sites, not all of which in turn have to be occupied by individuals. To identify the effects of viscous population structure, we compare two situations. In well-mixed populations, individuals interact with $n$ other individuals chosen at random from the entire population. In lattice-structured populations, sites are located on a lattice, with each individual occupying a site and interacting with individuals on $n$ neighboring sites. We used a square lattice with periodic boundary conditions, $30 \times 30$ sites, and the von Neumann neighborhood of $n = 4$ nearest neighbors.

Each individual $i$ possesses four adaptive traits ($c_i$, $c_{0,i}$, $p_{0,i}$, and $s_i$) that can all take continuous non-negative values. The cooperation level $c_i$ determines how much individual $i$ invests into cooperation with its neighbors: selfish individuals invest nothing or only a small amount, whereas cooperators invest a high amount. The punishment threshold $c_{0,i}$ determines the cooperation levels $c$ that individual $i$ deems sufficient or cooperative ($c > c_{0,i}$), as opposed to insufficient or selfish ($c < c_{0,i}$). Accordingly, selfish individuals with whom individual $i$ interacts are confronted with levels of punishment by individual $i$ that increase as their cooperation levels decrease. The punishment severity $p_{0,i}$ determines the punishment level individual $i$ metes out to individuals with a cooperation level of zero. The punishment strictness $s_i$ determines how sharply punishment by individual $i$ changes around $c_{0,i}$.

Each individual $i$ interacts with other individuals $j$ on $n$ neighboring sites in two steps: the interacting individuals cooperate according to their cooperation strategies and then punish according to their punishment strategies. The cooperation strategy of individual $i$ is given by its cooperation level $c_i$. For each investment $c_i$, individual $i$ pays the cooperation cost

$$C_c(c_i) = a_c c_i^{e_c}, \qquad (1)$$

with non-negative parameters $a_c$ and $e_c$. For $e_c < 1$ this cost function is decelerating, for $e_c = 1$ it is linear, and for $e_c > 1$ it is accelerating.

The punishment strategy of individual $i$ is given by its punishment reaction norm

$$p_i(c) = p_{0,i} \exp(-(c/c_{0,i})^{s_i}), \qquad (2a)$$

and depends on its punishment threshold $c_{0,i}$, punishment severity $p_{0,i}$, and punishment strictness $s_i$. This reaction norm describes the punishment level $p_i(c)$ with which individual $i$ responds to a cooperation level $c$. When punishment strictness $s_i$ is high, cooperation levels $c > c_{0,i}$ receive very little punishment, while cooperation levels $c < c_{0,i}$ elicit almost the maximal punishment level $p_{0,i}$. When punishment strictness $s_i$ is low, the punishment level still monotonically decreases as the cooperation level increases, but the transition to low punishment is shallower around $c_{0,i}$. For testing the robustness of our results, we also considered two alternative parameterizations of punishment reaction norms

$$p_i(c) = p_{0,i}(1 - c/c_{0,i})^{1/s_i} \quad \text{if } c < c_{0,i} \text{ and } p_i(c) = 0 \text{ otherwise}, \qquad (2b)$$

$$p_i(c) = p_{0,i}/[1 - \exp(-s_i) + \exp(s_i(c/c_{0,i} - 1))]. \qquad (2c)$$

In our model, punishment is costly. For each punishment level $p_i$, individual $i$ pays the punishment cost

$$C_p(p_i) = a_p p_i^{e_p}, \qquad (3)$$

with non-negative parameters $a_p$ and $e_p$. For $e_p < 1$ this cost function is decelerating, for $e_p = 1$ it is linear, and for $e_p > 1$ it is accelerating.

The birth rate of individual $i$

$$b_i = b_0 + \frac{1}{n}\sum_j c_j, \qquad (4a)$$

is given by the intrinsic birth rate $b_0$ increased by the average cooperative investment individual $i$ receives from is neighboring sites (the sums in Eqs. (4) extend over all individuals $j$ with whom individual $i$ interacts, and thus naturally exclude empty sites in the neighborhood of individual $i$). The resultant offspring is placed at a randomly chosen site with which individual $i$ is interacting, and is lost if that site is already occupied. Similarly, the death rate of individual $i$
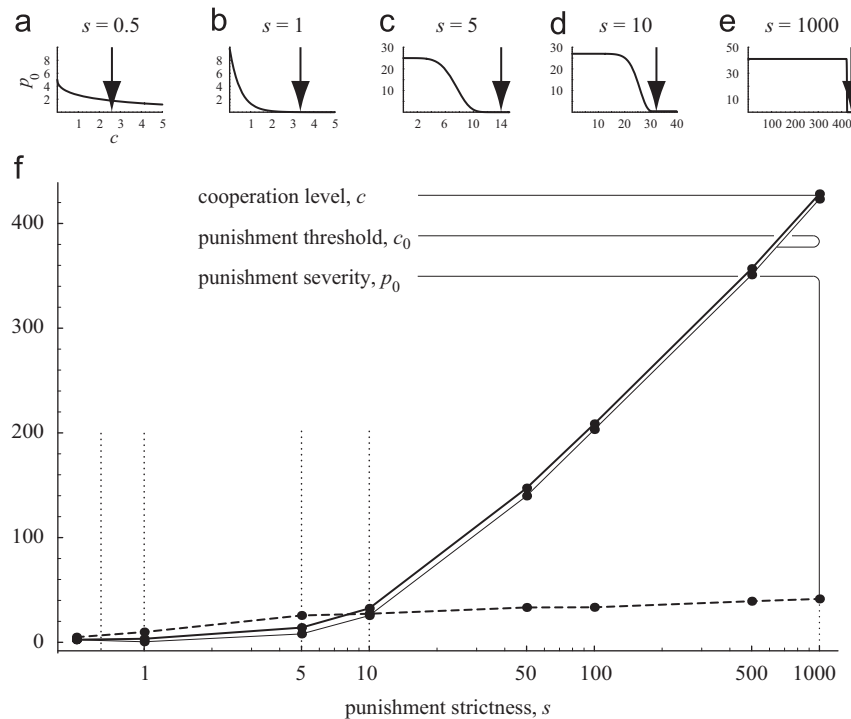
$$d_i = d_0 + \frac{1}{n}\sum_j [p_j(c_i) + C_p(p_i(c_j)) + C_c(c_i)] \qquad (4b)$$

is given by the intrinsic death rate $d_0$ increased by the average punishment individual $i$ receives and by the average costs for punishment and cooperation individual $i$ incurs.

Birth and death events occur asynchronously across the population and stochastically in time. After each such event, the waiting time until the next event is drawn from an exponential distribution with mean $1/E$ with $E = B + D$, where $B$ and $D$, respectively, are the current sums of all birth and death rates in the population. The event type is then chosen according to probabilities $B/E$ and $D/E$, and the individual $i$ undergoing the event is chosen according to probabilities $b_i/B$ or $d_i/D$.

When an offspring is born, its traits may be mutated relative to those of its parent. For each trait, a mutation occurs with probability $m$. Mutated trait values are normally distributed around the corresponding parental trait values, with standard deviations $\sigma_c$ for the traits $c$, $c_0$, and $p_0$, and with standard deviation $\sigma_s$ for the trait $s$. Mutated values of the traits $c$, $c_0$, $p_0$, and $s$ are constrained to minimal values 0, $10^{-5}$, 0, and 0, respectively. These boundaries are absorbing for $c$, $c_0$, and $p_0$, and reflective for $s$.

For testing the robustness of our results, we also considered errors in the implementation and perception of cooperation levels. With implementation errors, an implemented cooperation level differs from the actually intended cooperation level with a small error probability and with the difference being drawn from a normal distribution with a small standard deviation. With

**Fig. 1.** Joint evolution of cooperation level $c$, punishment threshold $c_0$, and punishment severity $p_0$, when punishment strictness $s$ is kept fixed. Panels (a)–(e) show the average evolved punishment reaction norms (continuous curves) and corresponding average evolved cooperation levels (vertical arrows) at time $t = 100,000$ for five different fixed values of punishment strictness $s$ (0.5, 1, 5, 10, and 1000). Panel (f) shows the average evolved values of $c$ (thick continuous curve), $c_0$ (thin continuous curve), and $p_0$ (dashed curve) as functions of $s$ (varying along the horizontal axis). All results are averaged over fifty model runs in the lattice-structured population. Runaway selection for cooperation and punishment accelerates with punishment strictness, leading to much elevated cooperation levels (for comparison: when punishment severity is kept fixed at $p_0 = 0$, the average cooperation level equilibrates at merely $c \approx 1.6$). The initial values of $c = 0$ and $p_0 = 0$ are chosen so as to highlight the bootstrapping of cooperation and punishment, i.e., their gradual and steady evolution in populations in which cooperation and punishment are entirely absent initially. The initial value of $c_0 = 10^{-5} > c$ means that all individuals are initially recognized as defectors. The initial frequency of empty sites is 50%. The punishment reaction norm is described by Eq. (2a). Other parameters: $b_0 = 2$, $d_0 = 1$, $a_c = 0.2$, $e_c = 1$, $a_p = 0.3$, $e_p = 0.5$, $m = 0.01$, and $\sigma_c = 1$.

perception errors, a perceived cooperation level differs from the actually implemented cooperation level analogously.
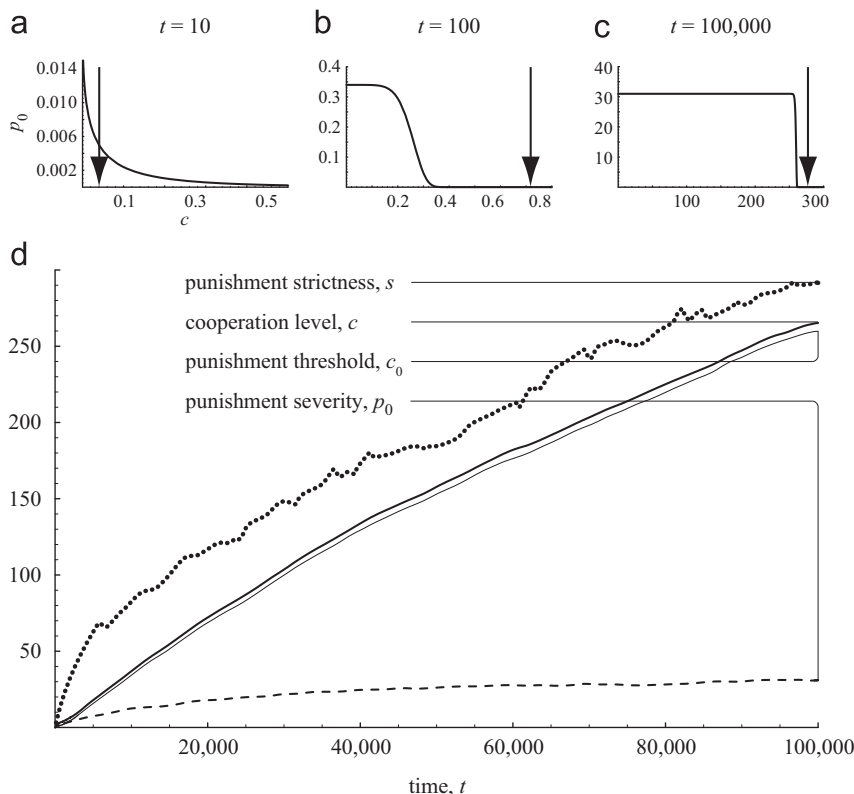
## 3. Results

Fig. 1 shows how our model leads to runaway selection for costly cooperation and punishment in lattice-structured populations. Here punishment strictness $s$ is not yet freely evolving, but instead is kept fixed at one and the same value for all individuals in the population. Evolution starts in the absence of any cooperation ($c = 0$) and of any punishment ($p_0 = 0$). All individuals are initially recognized as defectors ($c_0 = 10^{-5} > c$). In general, runaway selection among quantitative traits occurs when continual feedback between selection pressures and resultant evolutionary changes in the traits gradually leads to ever more extreme trait values. In our model, runaway selection occurs among the cooperation level $c$, the punishment threshold $c_0$, and the punishment severity $p_0$, which are all increasing concomitantly. We see that the larger $s$ is chosen, i.e., the stricter individuals apply their punishment threshold $c_0$, the faster these three traits evolve towards higher values. The population's average cooperation level $c$ always evolves to be slightly larger than the average punishment threshold $c_0$, so that most individuals are recognized as cooperators by most other individuals. Cooperation levels are driven up by evolutionary increases in punishment thresholds and vice versa. In other words, as the population evolves to become increasingly cooperative, the social demands on individuals to be recognized as cooperators rise

concomitantly. Also the punishment severity $p_0$ increases with the punishment strictness $s$. The speed of runaway selection thus increases with punishment strictness. Hence, stricter punishment indirectly favors both more severe punishment and higher cooperation levels.

Fig. 2 shows what happens when punishment strictness $s$ is allowed to evolve together with the three other adaptive traits $c$, $c_0$, and $p_0$. Again, evolution starts in the absence of any cooperation and of any punishment. In addition, individuals are assumed to be initially undiscriminating ($s = 0$). When the evolution of $s$ is sufficiently fast (i.e., when $\sigma_s$ is sufficiently large compared to $\sigma_c$), punishment strictness rises together with all other adaptive traits, resulting in a cooperative regime with strict-and-severe punishment. As in Fig. 1, the social requirements for avoiding punishment escalate with increasing cooperation. By contrast, when evolution of $s$ starts out from 0 but is too slow, punishment strictness remains low. Individuals thus continue to be undiscriminating, and runaway selection for cooperation and punishment cannot occur (results not shown). However, even when evolution of $s$ is slow, a sufficiently high initial value of $s$ reinstates the phenomenon of runaway selection, in line with the results already documented in Fig. 1.

Fig. 3 shows a systematic evaluation of the consequences of cooperation costs and punishment costs for the joint evolution of cooperation and punishment. Without punishment (i.e., for $p_0$ fixed at 0), cooperation evolves only when cooperation costs are sufficiently decelerating (Fig. 3a). Even then, resultant cooperation levels remain relatively low. Evolving punishment, by contrast, can lead to much higher levels of cooperation. This occurs when

**Fig. 2.** Joint evolution of cooperation level $c$, punishment threshold $c_0$, punishment severity $p_0$, and punishment strictness s. Panels (a)–(c) show the average evolved punishment reaction norms (continuous curves) and corresponding average cooperation levels (vertical arrows) at times $t = 10$, 100, and 100,000. Panel (d) shows the evolutionary dynamics of $c$ (thick continuous curve), $c_0$ (thin continuous curve), $p_0$ (dashed curve), and $s$ (dotted curve). The initial value of $s = 0.01$ implies an essentially flat reaction norm. Other parameters and settings are as in Fig. 1, with the addition of $\sigma_s = 10$.

punishment costs are decelerating or linear and cooperation costs are roughly linear (Fig. 3b). A look at the three traits determining the punishment strategy (Figs. 3c–e) confirms that these high levels of cooperation are enabled by the evolution of strict-and-severe punishment: the average punishment threshold (Fig. 3c) is again just slightly lower than the average cooperation level (Fig. 3b), the average punishment severity is high (Fig. 3d), and the average punishment strictness is also high (Fig. 3e).

We can categorize and understand these outcomes in terms of four cost scenarios. First, when cooperation is too cheap (i.e., cooperation costs are decelerating and $e_c$ is lower than about 0.5), the population's lattice structure alone is sufficient for promoting cooperation, so that costly punishment is not favored. Second, when cooperation is too expensive (i.e., cooperation costs are accelerating and $e_c$ is higher than about 1.25), cooperation evolution is hindered by these costs, independently of the costs of punishment. Third, when punishment is too expensive (i.e., punishment costs are accelerating and $e_p$ is higher than about 1.25), punishment evolution is hindered by these costs and no enhanced cooperation can thus occur. Fourth, when punishment is not too expensive (i.e., punishment costs are linear or decelerating so that $e_p$ is lower than about 1.25) and cooperation is neither too cheap nor too expensive (i.e., cooperation costs are roughly linear so that $e_c$ lies between about 0.5 and 1.25), runaway selection for cooperation and punishment occurs and results in greatly enhanced cooperation.
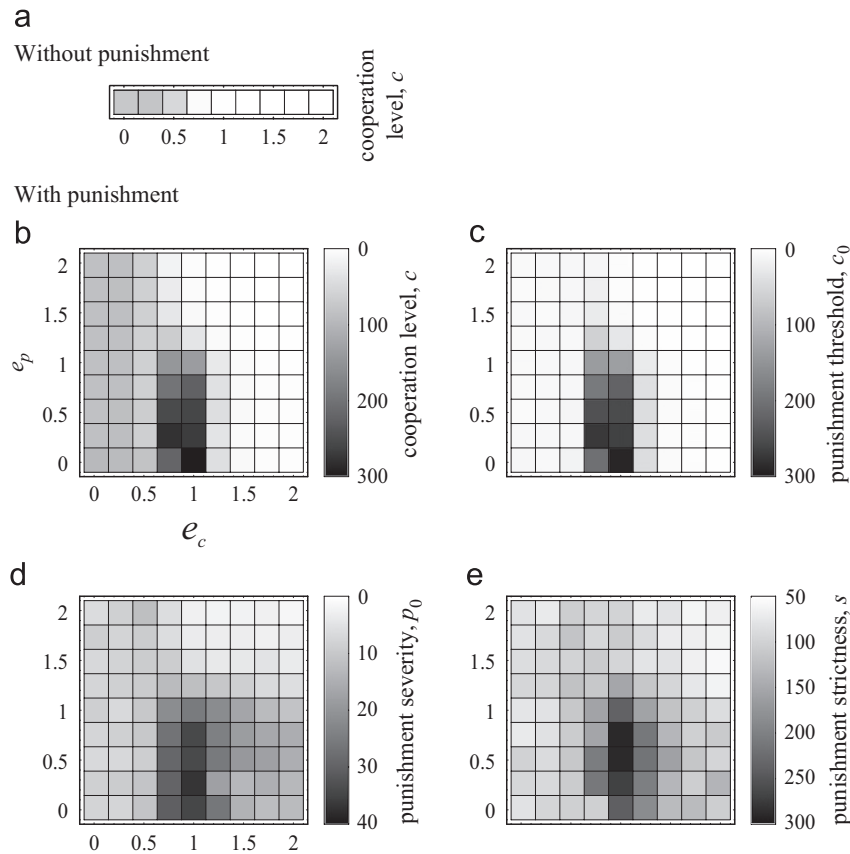
To test the robustness of our results, we changed the intrinsic birth and death rates, $b_0$ and $d_0$, without observing any qualitative differences. The patterns reported above also remain intact when we use the alternative parameterizations of punishment reaction norms in Eqs. (2b) and (2c), instead of the one in Eq. (2a). Also the introduction of implementation and perception errors did not lead

to any qualitative changes in the observed evolutionary dynamics. When increasing the mutation probability and the mutational standard deviations, we could confirm earlier results by Le Galliard et al. (2003) that showed how such changes in the mutation process facilitate the evolution of continuous cooperation strategies.
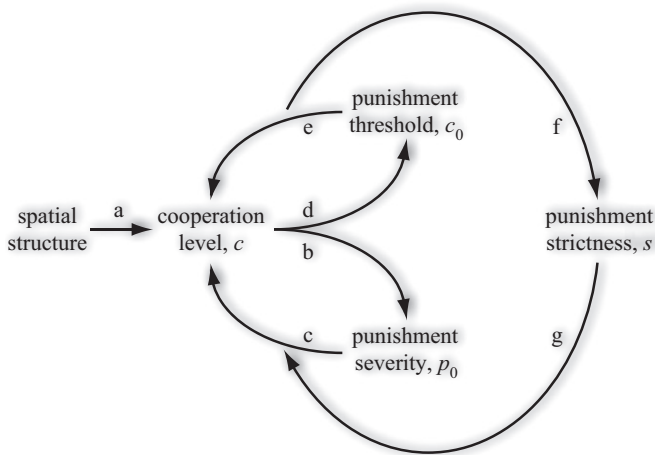
For well-mixed populations, the joint evolution of costly cooperation and punishment never occurs, as can be shown analytically (see Appendix A) and corroborated by individual-based simulations. This result can be understood intuitively: since punishing is costly to the punisher, and since in well-mixed populations this cost is the only selection pressure acting on punishment severity (see Appendix A), punishment—and, in its wake, cooperation—are invariably eliminated from well-mixed populations.

## 4. Discussion

Here we have shown that the joint and gradual evolution of cooperation and punishment can greatly promote cooperation levels in lattice-structured populations, even when cooperation and punishment are entirely absent initially. This promotion is driven by runaway selection, through which cooperation level, punishment threshold, and punishment severity rise concomitantly. The pace of the runaway process increases with punishment strictness. When punishment strictness is allowed to evolve, evolution often leads to strict-and-severe punishment accompanied by high cooperation levels. This process is again driven by runaway selection, now for all four traits. The enhancement of cooperation levels through the evolution of strict-and-severe punishment is largest when neither cooperation nor punishment

**Fig. 3.** Effects of cooperation and punishment costs on the joint evolution of cooperation level $c$, punishment threshold $c_0$, punishment severity $p_0$, and punishment strictness $s$. When the exponent $e_c$ ($e_p$) is small, equal to one, or large, costs for cooperation (punishment) are decelerating, linear, or accelerating. Decelerating (accelerating) costs imply that high levels of cooperation or punishment are relatively cheap (expensive). Panel (a) shows the average evolved cooperation level $c$ as a function of $e_c$ when punishment is absent (i.e., when $p_0$ is fixed at 0). Panels (b), (c), (d), and (e), respectively, show the average evolved values of $c$, $c_0$, $p_0$, and $s$ as functions of $e_c$ (varying along the horizontal axes) and $e_p$ (varying along the vertical axes). Other parameters and settings are as in Fig. 2.



**Fig. 4.** Schematic summary of positive feedbacks resulting in runaway selection for cooperation and strict-and-severe punishment.

are too costly and cooperation levels in the absence of punishment would be low. Our results explain the bootstrapping of cooperation and punishment, in the sense that defectors who rarely or only indiscriminately punish gradually and steadily evolve into cooperators who strictly and severely punish those they define as defectors.

The evolutionary mechanisms underlying these findings can be understood in intuitive terms. In general, any process of runaway

selection requires positive feedback between selection pressures and resultant evolutionary changes in one trait and selection pressures and resultant evolutionary changes in another trait. In our model, such mutual reinforcement can occur among all four evolving traits, as we have schematically summarized in Fig. 4. We start our explanation by recalling that lattice-structured populations enable the evolution of low levels of cooperation even in the absence of punishment (arrow a in Fig. 4). When punishment strictness is small but does not vanish completely, these cooperation levels favor increased punishment severity (arrow b). Under these conditions, punishment locally reduces the frequency of individuals with relatively low cooperation level, by differentially burdening them with a fitness disadvantage. Consequently, any region on the lattice in which punishment severity slightly differs from zero can expand into adjacent regions with vanishing punishment severity. Increased punishment severity then favors increased cooperation levels (arrow c), since these are advantageous when punishment reduces the exposure of more cooperative individuals to exploitation by less cooperative individuals. In turn, increased cooperation levels again favor increased punishment severity (arrow b), since this maintains the relative impact of punishment on fitness after cooperation levels have risen. Increased cooperation levels also favor increased punishment thresholds (arrow d), since this maintains the discrimination of individuals with relatively low cooperation levels after cooperation levels have risen. In turn, increased punishment thresholds favor increased cooperation levels (arrow e), since individuals must then cooperate more to escape punishment. Under these conditions, selection favors an

increase in punishment strictness (arrow f), since this enables a better targeting of punishment to individuals with relatively low cooperation levels. In turn, stricter punishment strengthens the already described selection pressures on cooperation level, punishment threshold, and punishment severity (arrow g), since stricter punishment selects for enhanced cooperation and tougher punishment.

These explanations help us to appreciate why runaway selection for cooperation and strict-and-severe punishment does not occur for all parameter values and initial conditions considered in our analysis. First, when the costs of cooperation or punishment are too high (upper and right regions in Figs. 3b–e), the selection pressures described above (arrows b–e in Fig. 4) are counteracted by those directly resulting from the costs, thus stalling the runaway process at low levels of cooperation and punishment. Second, when cooperation levels are high already in the absence of punishment (left regions in Figs. 3b–e), the relative advantages of punishment, and therefore the corresponding selection pressures on punishment (arrows b and d in Fig. 4), are low, thus stalling punishment evolution at low levels. Third, the initial selection pressure on punishment severity (arrow b in Fig. 4) occurs unless punishment is totally absent from the initial population. For the punishment reaction norms in Eq. (2b) the initial punishment threshold must thus exceed the initial cooperation level, since otherwise no punishment occurs at all. Fourth, for selection to favor stricter and severer punishment (arrows b, d, and f in Fig. 4), more cooperation has to result in less punishment, which implies that the punishment reaction norm must be a decreasing function. A vanishing punishment severity translates into a flat punishment reaction norm (Eqs. (2)), which prevents the runaway process from taking off. Conversely, this explains why increased punishment strictness accelerates the runaway process of the three other traits (Fig. 1) and why rapidly evolving punishment strictness facilitates the runaway process of all four traits (Fig. 2).

Our representation of cooperation and punishment strategies as continuous quantitative traits and the consideration of their gradual evolutionary dynamics play an important role for the findings reported here. In particular, the evolutionary mechanisms underlying the runaway process cause the steady and gradual adjustment of trait values driven by the subtle mutual reinforcement of selection pressures. In contrast, large sudden increases in punishment threshold or severity might not be selectively advantageous, since the resultant costs may outweigh the resultant benefits. Likewise, large sudden increases in cooperation levels are unlikely to be favored, since these would not be backed up by a corresponding orchestration of the punishment strategy. This highlights why cooperation games with continuous strategies and gradual trait evolution can reveal qualitative phenomena, such as the runaway selection for cooperation and strict-and-severe punishment reported here, that might be fundamentally obscured in corresponding games with discrete strategies.

Our results provide an evolutionary explanation for the widely observed appreciation of "strict but fair" punishment. This common cultural predisposition is an integral part of many moral systems and legal codes, and is often touted as a highly effective approach to education, reeducation, military discipline, and the preservation of public order. Strict-and-severe punishment is closely related to the "zero tolerance" approach to law enforcement, by which already small infractions of accepted rules are subjected to significant punishment. In our model, these ethical considerations have their counterpart in the emergence of high punishment strictness, elevated punishment severity, and of punishment thresholds finely tuned to majority behavior. In fact, our results presented in Figs. 1–3 make it clear that effective punishment must operate on shifting baselines, with the criterion

for punishment being continually refined as majority behavior evolves. Like in many other models of cooperation and punishment, these outcomes arise, gradually and naturally, from evolutionary dynamics solely driven by the selfish interests of individuals.

Based on these insights, we can revisit two conditions that could be perceived as limiting the bootstrapping of cooperation and punishment in our model. We had already explained above why runaway selection is hindered by vanishing initial punishment strictness, and, while punishment strictness is still low, by its low evolutionary rate. Notice that these observations only apply when punishment strictness is zero or very low initially. We can now question whether that would indeed be a realistic assumption. At least in humans, it seems fair to assume, instead, that innate or cultural circumstances are causing punishment strictness to start out from some intermediate level, even when punishment severity and punishment threshold start out from zero. Our results and explanations above make it clear that, under such conditions, runaway selection for cooperation and strict-and-severe punishment is greatly facilitated.

Here we have studied situations in which the punishment that individuals mete out simply depends on the cooperation levels of the individuals they interact with. Yet, punishment responses may be affected by many other factors. For example, breaking a social norm that is widely shared among members of a group may invite punishment (Gintis, 2000; Fehr and Fischbacher, 2004b), an effect that may be superimposed on the punishment responses considered here. Also emotions can influence punishment behavior, and may compel individuals to punish cheaters even when the cost of punishment exceeds that of being cheated (Frank, 1988; Xiao and Houser, 2005). Considering the effects of reputation or gossip on runaway selection for cooperation and punishment will also be of interest, since reducing an individual's reputation can serve as a cost-free means of punishment (Nakamaru and Kawata, 2004). Similarly, it will be worthwhile taking a closer look at conditions and mechanisms that can eventually stop the runaway process investigated here. This could involve cost functions that are decelerating for low investments and accelerating for high investments, diminishing fitness returns from received investments, or an explicit modeling of the availability of resources that individuals exchange when they cooperate or punish.

The evolutionary framework we have utilized here recognizes three levels of interlocking dynamics, ranging from the demographic dynamics of individuals in a population, to the behavioral dynamics of cooperation and punishment in the interactions between individuals, and to the psychological dynamics underlying the identification of cheaters. Naturally, psychological dynamics affect behavioral dynamics, which in turn affect demographic dynamics. Conversely, demographic dynamics affect behavioral and psychological dynamics by changing the selection pressures that cause adjustments in the traits governing behavior and psychology. Experimental tools and modeling approaches for studying such feedbacks have emerged over the past decades and are now increasingly applied to tackling questions in cooperation research (e.g., de Quervain et al., 2004; Enquist and Ghirlanda, 2005). We hope that the framework and results put forward here may further inspire and facilitate such studies. In a similar vein, our approach could be used to address questions raised by evolutionary psychologists who have challenged conjectured adaptive explanations of behavior and psychological predispositions regarding mate choice, emotion, cheater detection, and the ability to recognize spatial locations (e.g., Bawkow et al., 1992). While such explanations are often based on verbal and qualitative reasoning, the approach adopted here allows for formal and quantitative reasoning.

It is our hope that, from a methodological perspective, our evolutionary explanation of runaway selection for cooperation and strict-and-severe punishment might be no more than a start. We believe that, more in general, studies of cooperation have much to gain from investigating models with joint evolution of multiple continuous traits, explicit dynamics for demography and trait changes, and interpretation of traits in terms of reaction norms for psychological and behavioral processes.

## Acknowledgements

## Appendix A

In this appendix we show that cooperation and punishment cannot evolve in well-mixed populations. For this purpose we investigate the dynamics of a rare variant strategy with frequency $x' \approx 0$, cooperation level $c'$, and punishment reaction norm $p'$ in the population of a resident strategy with frequency $x$, cooperation level $c$, and punishment reaction norm $p$

$$\frac{1}{x'}\frac{dx'}{dt} = (b_0 + cx)(1-x) - \{d_0 + [p(c') + C_p(p'(c)) + C_c(c')]x\}.$$

We assume that the resident population is at its equilibrium frequency $0 \leqslant \hat{x} \leqslant 1$, so that $(b_0 + c\hat{x})(1-\hat{x}) = d_0 + [p(c) + C_p(p(c)) + C_c(c)]\hat{x}$, from which we obtain

$$\hat{x} = \frac{1}{2c}\left[\sqrt{l^2 + 4c(b_0 - d_0)} - l\right]$$

with $l = b_0 - c + p(c) + C_p(p(c)) + C_c(c)$. Denoting the variant's per capita growth rate or fitness $(dx'/dt)/x'$ by $f'$ (e.g., Metz et al., 1992), the selection pressures $g_c$, $g_{c_0}$, $g_{p_0}$, and $g_s$ on the resident's adaptive traits $c$, $c_0$, $p_0$, and $s$ are given by the derivatives $df'/dc'$, $df'/dc_0'$, $df'/dp_0'$, and $df'/ds'$ evaluated at $c' = c$ and $p' = p$ (e.g., Dieckmann and Law, 1996; Geritz et al., 1997). Using Eqs. (1), (2a), and (3), this gives

$$g_c = \hat{x}c^{-1}[sc^s c_0^{-s}p(c) - e_c C_c(c)],$$

$$g_{c_0} = -\hat{x}sc^s c_0^{-(s+1)}e_p C_p(p(c)),$$

$$g_{p_0} = -\hat{x}p_0^{-1}e_p C_p(p(c)),$$

$$g_s = \hat{x}c^s c_0^{-s}e_p C_p(p(c))\ln(c/c_0).$$

Since $g_{p_0}$ is negative, evolution will always diminish punishment severity $p_0$ in well-mixed populations. Once $p_0$ has evolved to 0, selection on $c_0$ and $s$ ceases: $C_p(0) = 0$ and thus $g_{c_0} = 0$ and $g_s = 0$. The selection pressure on $c$ is negative for $p_0 = 0$, $g_c = -\hat{x}c^{-1}e_c C_c(c)$, so that, driven by the cost of cooperation, the cooperation level $c$ will also evolve to 0.

## References

Axelrod, R., 1986. An evolutionary approach to norms. Am. Political Sci. Rev. 80, 1095–1111.

Axelrod, R., Hamilton, W.D., 1981. The evolution of cooperation. Science 211, 1390–1396.

Bawkow, J.H., Cosmides, L., Tooby, J., 1992. The Adapted Mind: Evolutionary Psychology and the Generation of Culture. Oxford University Press, Oxford.

Bowles, S., Gintis, H., 2004. The evolution of strong reciprocity: cooperation in heterogeneous populations. Theor. Popul. Biol. 65, 17–28.

Boyd, R., Richerson, P.J., 1992. Punishment allows the evolution of cooperation (or anything else) in sizable groups. Ethology Sociobiology 13, 171–195.

Boyd, R., Gintis, H., Bowles, S., Richerson, P.J., 2003. The evolution of altruistic punishment. Proc. Natl. Acad. Sci. USA 100, 3531–3535.

Brandt, H., Sigmund, K., 2004. The logic of reprobation: assessment and action rules for indirect reciprocation. J. Theor. Biol. 231, 475–486.

Brandt, H., Hauert, C., Sigmund, K., 2003. Punishment and reputation in spatial public goods games. Proc. R. Soc. London B 270, 1099–1104.

Brandt, H., Hauert, C., Sigmund, K., 2006. Punishing and abstaining for public goods. Proc. Natl. Acad. Sci. USA 103, 495–497.

Clutton-Brock, T.H., Parker, G.A., 1995. Punishment in animal societies. Nature 373, 209–216.

de Quervain, D.J.F., Fischbacher, U., Treyer, V., Schellhammer, M., Schnyder, U., Buck, A., Fehr, E., 2004. The neural basis of altruistic punishment. Science 305, 1254–1258.

Dieckmann, U., Law, R., 1996. The dynamical theory of coevolution: a derivation from stochastic ecological processes. J. Math. Biol. 34, 579–612.

Doebeli, M., Knowlton, N., 1998. The evolution of interspecific mutualisms. Proc. Natl. Acad. Sci. USA 95, 8676–8680.

Doebeli, M., Hauert, C., Killingback, T., 2004. The evolutionary origin of cooperators and defectors. Science 306, 859–862.

Eldakar, O.T., Wilson, D.S., 2008. Selfishness as second-order altruism. Proc. Natl. Acad. Sci. USA 105, 6982–6986.

Eldakar, O.T., Farrell, D.L., Wilson, D.S., 2007. Selfish punishment: altruism can be maintained by competition among cheaters. J. Theor. Biol. 249, 198–205.

Enquist, M., Ghirlanda, S., 2005. Neural Networks and Animal Behavior. Princeton University Press, Princeton.

Fehr, E., Fischbacher, U., 2004a. Third-party punishment and social norms. Evol. Hum. Behav. 25, 63–87.

Fehr, E., Fischbacher, U., 2004b. Social norms and human cooperation. Trends Cognitive Sci. 8, 185–190.

Fehr, E., Gächter, S., 2002. Altruistic punishment in humans. Nature 415, 137–140.

Fehr, E., Rockenbach, B., 2003. Detrimental effects of sanctions on human altruism. Nature 422, 137–140.

Fowler, J.H., 2005. Altruistic punishment and the origin of cooperation. Proc. Natl. Acad. Sci. USA 19, 7047–7049.

Frank, R.H., 1988. Passions within Reason. W.W. Norton, New York.

Gardner, A., West, S.A., 2004. Cooperation and punishment, especially in humans. Am. Nat. 164, 753–764.

Geritz, S.A.H., Metz, J.A.J., Kisdi, É., Meszéna, G., 1997. Dynamics of adaptation and evolutionary branching. Phys. Rev. Lett. 78, 2024–2027.

Gintis, H., 2000. Strong reciprocity and human sociality. J. Theor. Biol. 206, 169–179.

Hamilton, W.D., 1964. The genetical evolution of social behavior. I. J. Theor. Biol. 7, 1–52.

Hauert, C., Traulsen, A., Brandt, H., Nowak, M.A., Sigmund, K., 2007. Via freedom to coercion: the emergence of costly punishment. Science 316, 1905–1907.

Henrich, J., Boyd, R., 2001. Why people punish defectors? Weak conformist transmission can stabilize costly enforcement of norms in cooperative dilemmas. J. Theor. Biol. 208, 79–89.

Henrich, J., McElreath, R., Barr, A., Ensminger, J., Barrett, C., Bolyanatz, A., Cardenas, J.C., Gurven, M., Gwako, E., Henrich, N., Lesorogol, C., Marlowe, F., Tracer, D., Ziker, J., 2006. Costly punishment across human societies. Science 312, 1767–1770.

Killingback, T., Doebeli, M., 2002. The continuous Prisoner's Dilemma and the evolution of cooperation through reciprocal altruism with variable investment. Am. Nat. 160, 421–438.

Killingback, T., Doebeli, M., Knowlton, N., 1999. Variable investment, the continuous Prisoner's Dilemma, and the origin of cooperation. Proc. R. Soc. London B 266, 1723–1728.

Le Galliard, J., Ferrière, R., Dieckmann, U., 2003. The adaptive dynamics of altruism in spatially heterogeneous populations. Evolution 57, 1–17.

Le Galliard, J., Ferrière, R., Dieckmann, U., 2005. Adaptive evolution of social traits: origin, trajectories, and correlations of altruism and mobility. Am. Nat. 165, 206–224.

Lehmann, L., Keller, L., West, S., Roze, D., 2007a. Group selection and kin selection: two concepts but one process. Proc. Natl. Acad. Sci. USA 104, 6736–6739.

Lehmann, L., Rousset, F., Roze, D., Keller, L., 2007b. Strong reciprocity or strong ferocity? A population genetic view of the evolution of altruistic punishment. Am. Nat. 170, 21–36.

Leimar, O., Hammerstein, P., 2001. Evolution of cooperation through indirect reciprocity. Proc. Roy. Soc. London B 268, 745–753.

Matsuda, H., 1987. Condition for the evolution of altruism. In: Ito, Y., Brown, J., Kikkawa, J. (Eds.), Animal Societies: Theories and Facts. Scientific Society Press, Tokyo, pp. 56–79.

Metz, J.A.J., Nisbet, R.M., Geritz, S.A.H., 1992. How should we define "fitness" for general ecological scenarios? Trends Ecol. Evol. 7, 198–202.

Nakamaru, M., Iwasa, Y., 2005. The evolution of altruism by costly punishment in lattice-structured populations: score-dependent viability versus score-dependent fertility. Evol. Ecol. Res. 7, 853–870.

Nakamaru, M., Iwasa, Y., 2006. The coevolution of altruism and punishment: role of the selfish punisher. J. Theor. Biol. 240, 475–488.

Nakamaru, M., Kawata, M., 2004. Evolution of rumors that discriminate lying defectors. Evol. Ecol. Res. 6, 261–283.

Nakamaru, M., Matsuda, H., Iwasa, Y., 1997. The evolution of cooperation in a lattice-structured population. J. Theor. Biol. 184, 65–81.

Nakamaru, M., Nogami, H., Iwasa, Y., 1998. Score-dependent fertility model for the evolution of cooperation in a lattice. J. Theor. Biol. 194, 101–124.

Nowak, M.A., 2006. Five rules for the evolution of cooperation. Science 314, 1560–1563.

Nowak, M.A., May, R.M., 1992. Evolutionary games and spatial chaos. Nature 359, 826–829.

Nowak, M.A., Sigmund, K., 1998. Evolution of indirect reciprocity by image scoring. Nature 393, 573–577.

Ohtsuki, H., Iwasa, Y., 2004. How should we define goodness?—reputation dynamics in indirect reciprocity. J. Theor. Biol. 231, 107–120.

Ohtsuki, H., Hauert, C., Lieberman, E., Nowak, M.A., 2006. Simple rule for the evolution of cooperation on graphs and social networks. Nature 441, 502–505.

Panchanathan, K., Boyd, R., 2003. A tale of two defectors: the importance of standing for evolution of indirect reciprocity. J. Theor. Biol. 224, 115–126.

Roberts, G., Sherratt, T.N., 1998. Development of cooperative relationships through increasing investment. Nature 394, 175–179.

Rockenbach, B., Milinski, M., 2006. The efficient interaction of indirect reciprocity and costly punishment. Nature 444, 718–723.

Shinada, M., Yamagishi, T., Ohmura, Y., 2004. False friends are worse than bitter enemies: "altruistic" punishment of in-group members. Evol. Hum. Behav. 25, 379–393.

Sigmund, K., 2007. Punish or perish? Retaliation and collaboration among humans. Trends Ecol. Evol. 22, 593–600.

Sigmund, K., Hauert, C., Nowak, M.A., 2001. Reward and punishment. Proc. Natl. Acad. Sci. USA 98, 10757–10762.

Sober, E., Wilson, D.S., 1998. Unto Others: The Evolution and Psychology of Unselfish Behavior. Harvard University Press, Cambridge.

Takahashi, N., Mashima, R., 2006. The importance of subjectivity in perceptual errors on the emergence of indirect reciprocity. J. Theor. Biol. 243, 418–436.

Wahl, L.M., Nowak, M.A., 1999a. The continuous Prisoner's Dilemma: I. Linear reactive strategies. J. Theor. Biol. 200, 307–321.

Wahl, L.M., Nowak, M.A., 1999b. The continuous Prisoner's Dilemma: II. Linear reactive strategies with noise. J. Theor. Biol. 200, 323–338.

Xiao, E., Houser, D., 2005. Emotion express in human punishment behavior. Proc. Natl. Acad. Sci. USA 102, 7398–7401.

Yamagishi, T., 1986. The provision of a sanctioning system as a public good. J. Pers. Soc. Psychol. 51, 110–116.