Original Article

## Collective Action and the Detrimental Side of Punishment

Shade T. Shutters, School of Sustainability and Center for Social Dynamics and Complexity, Arizona State University, Tempe, Arizona, USA. Email: shade.shutters@asu.edu (Corresponding author).

**Abstract:** Cooperative behavior is the subject of intense study in a wide range of scientific fields, yet its evolutionary origins remain largely unexplained. A leading explanation of cooperation is the mechanism of altruistic punishment, where individuals pay to punish others but receive no material benefit in return. Experiments have shown such punishment can induce cooperative outcomes in social dilemmas, though sometimes at the cost of reduced social welfare. However, experiments typically examine the effects of punishing low contributors without allowing others in the environment to respond. Thus, the full ramifications of punishment may not be well understood. Here, I use evolutionary simulations of agents playing a continuous prisoners dilemma to study behavior subsequent to an act of punishment, and how that subsequent behavior affects the efficiency of payoffs. Different network configurations are used to better understand the relative effects of social structure and individual strategies. Results show that when agents can either retaliate against their punisher, or punish those who ignore cheaters, the cooperative effects of punishment are reduced or eliminated. The magnitude of this effect is dependent on the density of the network in which the population is embedded. Overall, results suggest that a better understanding of the aftereffects of punishment is needed to assess the relationship between punishment and cooperative outcomes.

**Keywords:** strong reciprocity, cooperation, altruism, retaliation, simulation, network, social behavior.

--------------------------------------------------------------------------------

## Introduction

Explaining the evolution of cooperation is one of the greatest unanswered questions facing evolutionary biologists today (Maynard Smith and Szathmáry, 1997; West, Griffin, and Gardner, 2007). Cooperation is instrumental in maintaining human social institutions (Ostrom, Walker, and Gardner, 1992) and is required among nations to effectively address global-scale problems (Kaul and Mendoza, 2003; Sandler, 1997). Thus, understanding the mechanisms that result in cooperation is important to both scientists and policy makers.

Yet, despite its fundamental importance, the evolution of cooperative behavior remains largely unexplained.

Several mechanisms have been previously suggested to explain the evolution of cooperation, including kin selection (Hamilton, 1964), multi-level selection (Fletcher and Zwick, 2004; Goodnight, 2005; Reeve and Hölldobler, 2007; Wilson and Wilson, 2007; Wilson and Hölldobler, 2005), direct reciprocity (Axelrod and Hamilton, 1981; Trivers, 1971), indirect reciprocity (Boyd and Richerson, 1989; Leimar and Hammerstein, 2001; Nowak and Sigmund, 2005), and tag-mediated altruism (Axelrod, Hammond, and Grafen, 2004; Riolo, Cohen, and Axelrod, 2001; Spector and Klein, 2006). While these mechanisms explain some instances of cooperation, they generally apply to limited cases or special circumstances such as genetic relatedness or long-term relationships between individuals. The search for a more broadly applicable explanation has increasingly focused on altruistic punishment, where individuals incur a cost to punish others without receiving any material benefit in return.

Punishment is ubiquitous among social organisms and wherever cooperating individuals have an incentive to cheat or free-ride, punishment behavior usually exists as a deterrent (Frank, 1995). This includes toxin release in colonial bacteria that affects only non-cooperators (Travisano and Velicer, 2004), the destruction of eggs laid by workers in social insect colonies (Foster and Ratnieks, 2001), and enforcement of dominance and mating hierarchies in non-human mammals (Clutton-Brock and Parker, 1995). Even the process of cellular meiosis can be viewed as a form of policing selfish genes (Michod, 1996). In humans, punishment and policing are common across many diverse societies and cultural groups (Marlowe et al., 2008) and are prevalent in local-scale management of common property (Coleman and Steed, 2009). Policy makers view punishment institutions as key to resolving social conflict both at local scales, in governance of common pool resources (Dietz, Ostrom, and Stern, 2003; Ostrom et al., 1992), and at global scales where it is considered a prerequisite for successful international agreements (Barrett, 2003). Laboratory and simulation experiments generally support the idea that altruistic punishment can lead to the provisioning of public goods (Boyd, Gintis, Bowles, and Richerson, 2003; Fehr and Gächter, 2000; Gürerk, Irlenbusch, and Rockenbach, 2006; Ostrom et al., 1992; Shutters, 2012), though others have demonstrated exceptions to this finding (Cinyabuguma, Page, and Putterman, 2006; Fehr and Rockenbach, 2003; Herrmann, Thöni, and Gächter, 2008)

While it is fine to propose that altruistic punishment is a mechanism leading to the evolution of cooperation, this only shifts the underlying question from "why should an individual cooperate?" to "why should an individual altruistically punish?" As research begins to focus on the latter question, cultural group selection (Hagen and Hammerstein, 2006; Richerson and Boyd, 2005) and the feedbacks of social structures (Shutters, 2012) have been recently suggested as mechanisms leading to the evolution of altruistic punishment.

What has not been adequately addressed is how punishment affects the efficiency of cooperation, a measure of the net increase in payoffs that result when punishment is used to induce cooperation (Nikiforakis, 2008; Sefton, Shupp, and Walker, 2007). Even if punishment induces a society to cooperate there are costs associated with punishing and

being punished that reduce the overall gains from cooperation, and these must be accounted for when discussing the efficacy of punishment as a cooperative mechanism.

Previous experiments using punishment show that its use can increase contributions to a public even though total payoffs decrease compared to a population comprised only of defectors (Fehr and Gächter, 2000; Ostrom et al., 1992). This negative affect on efficiency has been demonstrated when interactions are not repeated sufficiently, though increasing the number of repeat interactions eventually led to positive gains in total payoffs (Gächter, Renner, and Sefton, 2008; Gürerk et al., 2006). Thus, it remains unclear in a world of increasingly frequent one-shot interactions how punishment used to induce cooperation might affect total payoffs.

Understanding how punishment-induced cooperation affects payoff efficiency is especially important when considering the aftermath of punishment. Experiments with punishment typically include only a round of game play and a round in which agents can punish cheaters. These experiments ignore the fact that punishment in real-world situations usually elicits further responses of some type.

Thus, the purpose of the current study is not to support or refute mechanisms that may lead to the evolution of altruistic punishment. Instead, its purpose is to understand how the efficiency of punishment-induced cooperation is affected when a more realistic range of behavior is allowed to take place. In this study, the specific behaviors of retaliation and second-order punishment are allowed in a simulated society and their effects on the efficiency of cooperation are examined.

*Second-order punishment*

Sanctioning and policing institutions often exist in human societies to deter cheating in the provisioning of public goods. But a paradox arises, known as the second-order free-rider problem, regarding what motivates those who are supposed to punish cheaters (Hodgson, 2009; Sigmund, De Silva, Traulsen, and Hauert, 2010). Without deterrents and/or incentives, enforcement agents are expected to avoid the costs and risks of punishing and to simply ignore cheaters. These agents that avoid their policing duties have an evolutionary advantage over those that do punish (Dreber, Rand, Fudenberg, and Nowak, 2008) and the mechanism of second-order punishment often exists as a deterrent against policing agents that do not do their jobs. Second-order punishment occurs when an agent declines to punish cheaters when given the chance, and is itself punished as a result (Boyd and Richerson, 1992). Even though such individuals may otherwise cooperate and contribute substantially to a public good, they are punished because they take no action against cheaters.

But what is the effect on payoff efficiencies when agents are seemingly coerced into punishing cheaters? One may reason that, since punishing cheaters induces public good contributions, punishing those that ignore cheating will only further enhance public good contributions. On the other hand, laboratory experiments with human participants have demonstrated the opposite, showing that sanctioning otherwise cooperative agents because they ignore cheating can inhibit the emergence of cooperation (Denant-Boemont, Masclet, and Noussair, 2007). This leads to the first question addressed in this study:

Q1) When agents may altruistically punish others that permit cheating, how is efficiency of cooperation affected?

*Retaliation*

Another under-addressed behavior that often co-occurs with punishment is retaliation. Research has shown that humans and other animals are not indifferent to being punished and often retaliate at a cost to both themselves and their punisher (Clutton-Brock and Parker, 1995; Molm, 1994). The prospect of suffering retaliation can deter agents from punishing free-riders (Nikiforakis and Engelmann, 2011) and ultimately negate the cooperative effects of punishment (Nikiforakis, 2008). This consequence is frequently overlooked in studies of punishment-induced cooperation (Denant-Boemont et al., 2007), which typically allow only punishment of cheaters and do not allow a response from the punished party. Thus the third simulation allows the ability to retaliate when punished and seeks to answer the question:

Q2) When a punished agent may retaliate against its punisher, how are aggregate levels of cooperation affected compared to simulations without retaliation?

*Social welfare*

In both treatments, second-order punishment and retaliation, the focus of this study is not simply on how contributions to a public good are affected, but on how a population's overall payoffs are affected. Increased contributions to a public good are typically assumed to be due to cooperative behavior but it may also be that contributions increase because of coercion. This is an important distinction that becomes clearer when analyzing how a given treatment affects total net payoffs or payoff efficiency. This study draws a distinction between cooperation, increasing both contributions and payoffs, and coercion, increasing contributions at the expense of lower payoffs. Thus, this study also seeks to answer the question:

Q3) If either retaliation or second-order punishment induces higher levels of cooperation does it also increase aggregate payoffs?

*Population structure*

Research has demonstrated that populations embedded in spatially explicit grids can evolve different aggregate attributes than non-structured populations (Boyd and Richerson, 2002; Killingback and Doebeli, 1996; Killingback and Studer, 2001; Nowak and May, 1992; Page, Nowak, and Sigmund, 2000). More importantly, several studies show that network structure plays a critical role in the evolution of aggregate behavior such as cooperation (Chen, Fu, and Wang, 2007; Chwe, 1999; Gould, 1993; Huang, Wang, Xu, and Wang, 2008; Hui, Xu, and Zheng, 2007; Ifti, Killingback, and Doebelic, 2004; Ohtsuki, Hauert, Lieberman, and Nowak, 2006; Santos, Rodrigues, and Pacheco, 2006), especially when those networks are dynamic and coevolving with the agents they govern (Hales, 2005; Santos, Pacheco, and Lenaerts, 2006; Shutters and Cutts, 2008; Takács, Janky, and Flache, 2008).

Therefore, it is critical to understand not only how second-order punishment and retaliation affect the efficiency of cooperation but also how differences in population structure influence outcomes. This study examines the role of social structure by analyzing evolutionary outcomes both with and without structured societies.

## Materials and Methods

To test the questions outlined above, various punishment options were incorporated into evolutionary simulations of the continuous prisoner's dilemma. Social simulations, including agent-based models, individual-based models, and other evolutionary computational techniques, offer unique insights into dynamic behavior (North, 2005) such as the relationship between individual behavior and emergent properties at higher scales (Anderies, 2002; Harrison and Singer, 2006), that are typically not provided by formal models of social systems (Harrison and Singer, 2006; Sawyer, 2005). Social simulations also allows careful control over factors that may confound empirical studies such as emotion, reputation, visual cues, anonymity, or cultural influences (Cederman, 2001), while probing vast expanses of evolutionary space that would be impractical in laboratory settings due to cost or time constraints. It should be noted that social simulations are generally designed as a complement to laboratory experiments and cast studies, not as a replacement.

To understand the effects of social structure, which are known to significantly influence results of social simulations (Santos, Rodrigues, and Pacheco, 2006), simulations were conducted both with and without social structure. When added, social structure consisted of regular networks of varying density.

*The continuous prisoner's dilemma (CPD)*

In the standard prisoner's dilemma players are limited to two choices - cooperate or defect. Here, that requirement is relaxed and players select a level of cooperation on the continuum between full cooperation and full defection. This presents an arguably more realistic picture of choices facing those in social dilemmas (Killingback and Doebeli, 2002; Sandler, 1999) and is known as the continuous prisoner's dilemma (CPD).

In a CPD game $i$ and $j$ are each given an endowment standardized here to one unit. From this they independently and simultaneously contribute a portion $x \in [0,1]$ to a public good pool, while keeping the remainder, so that $x = 1$ represents full cooperation and $x = 0$ full defection (Deng and Chu, 2011; Schofield, 1977). For any given contribution by $j$, $i$'s payoff is maximized when $x_i = 0$. This is the expected rational choice or Nash equilibrium of the CPD. The dilemma arises, however, because total social welfare, measured as total net payoffs, is maximized when both individuals cooperate fully and $x_i = x_j = 1$.

*Social Structure*

At the beginning of each simulation, a specific network is generated that structures the population and determines the allowable interactions between agents. All networks are non-directed, unweighted, and static.

To understand the effects of social structure on outcomes, a number of regular

networks are used. Often represented as lattice structure, are those in which all nodes have the same degree $d$ (or number of neighbors) and are arranged in a regular repeating pattern. In addition, these networks are torroidal, meaning that they have no edges but instead loop around onto themselves such as the surface of a sphere. Two regular networks used commonly in simulations, including this study, are the von Neumann network ($d = 4$) and the Moore network ($d = 8$). Hexagonal networks ($d = 6$) are also used as well as one-dimensional rings known as linear networks ($d = 2$). Though regular networks bear little resemblance to interaction patterns in real-world social systems, their use in simulation studies reduces confounding effects of social structure because they have no variance in degree, no edge effects, and uniform distances among individuals in a population. When used in this study, regular networks are referenced throughout this paper by their degree $d$.

In contrast to structured societies, complete networks are used in this study to understand how the absence of social structure affects outcomes. Complete networks are those in which every node is linked to every other node in the population. Though technically a regular network with $d = N - 1$, where $N$ is the population size, an agent in a complete network has equal probability of interacting with any other agent. Thus, complete networks are analogous to homogeneous, well-mixed systems that have no social structure. Throughout this study, simulations using complete networks are synonymous with unstructured populations.

*Base game play*

In the base game, agents play the CPD followed by a single round of punishment. A single simulation run initiates with creation of a social network. Each node is occupied by a single agent $i$ consisting of strategy ($x_i$, $t_i$, $c_i$, $s_i$) where $x_i$ = the contribution $i$ makes to the public good in the CPD, $t_i$ = the contribution below which the agent will punish another agent in a game being observed by $i$, $c_i$ = how much $i$ spends to punish an observed agent whose contribution is too low, and $s_i$ = the amount $i$ spends to retaliate when it has been punished (in simulations that allow retaliation). In other words, $t_i$ determines *if* agent $i$ will punish and $c_i$ determines *how much* agent $i$ will punish. Each strategy component $x_i$, $t_i$, $c_i$, $s_i$ ∈ [0,1] and is generated randomly from a uniform distribution at the beginning of each simulation. To control for other factors that might contribute to the maintenance of cooperation, such as history or reputation, agents have no memory of prior interactions or agents. Every game is effectively one-shot and anonymous.

During a single CPD game an agent $i$ initiates the encounter by randomly selecting $j$ from its neighborhood, which consists of all nodes one link away from $i$ in the given network type. Agents are given their endowment of one unit from which each simultaneously contributes a portion to the public good pool. Payoffs are then calculated as in Table 1. The initiating player $i$ then randomly selects a second neighbor $k$, who is tasked with observing and evaluating $i$'s contribution. If $k$ judges the contribution to be too low ($x_i < t_k$), $k$ pays $c_k$ to punish $i$ by the amount $c_kM$, where $M$ is the relative strength of punishment referred to here as the punishment multiplier. During a single generation of the simulation each agent initiates three CPD games and, on average, acts as an observer (and possible punisher) three times. A single simulation run execute for 10,000 generations.

Each generation consists of three routines - game play, observation and punishment

(including retaliation and punishment of non-punishers), and selection and reproduction. During each routine an agent interacts only with its immediate neighbors as defined by the network type and all interactions take place in parallel. For each agent, $p$ represents the net payoffs (benefits − costs) an agent earns during a generation. At the start of a new generation $p = 0$. It is increased by the amount earned in each CPD but is decreased when the agent is punished by other agents or when the agent pays to retaliate or punish someone else.

**Table 1.** Payoffs $p$ in a CPD with:

(a) players $i$ and $j$ and punisher $k$ observing and possibly punishing $i$

| Payoff | $x_i \geq t_k$   ($k$ does not punish $i$) | $x_i < t_k$   ($k$ punishes $i$) |
|---|---|---|
| $p_i$ | $1 - x_i + r(x_i + x_j)/2$ | $1 - x_i + r(x_i + x_j)/2 - c_k M$ |
| $p_j$ | $1 - x_j + r(x_i + x_j)/2$ | $1 - x_j + r(x_i + x_j)/2$ |
| $p_k$ | $0$ | $- c_k$ |

(b) punisher $k$ being observed and possibly punished by second-order punisher $l$

| Payoff | $t_l < t_k$   ($l$ does not punish $k$) | $t_l \geq t_k$   ($l$ punishes $k$) |
|---|---|---|
| $p_k$ | $0$ | $- c_l M$ |
| $p_l$ | $0$ | $- c_l$ |

(c) player $i$ retaliating when punished by $k$

| Payoff | $x_i \geq t_k$   ($k$ does not punish $i$) | $t_l \geq t_k$   ($k$ punishes, $i$ retaliates) |
|---|---|---|
| $p_i$ | $1 - x_i + r(x_i + x_j)/2$ | $1 - x_i + r(x_i + x_j)/2 - c_k M - x_i$ |
| $p_k$ | $0$ | $- c_k - x_i M$ |

Following game play and punishment, agents compete with one another in the reproduction routine for the ability to pass offspring to the next generation. During this routine each agent $i$ randomly selects a neighbor $j$ with which to compare respective payoffs accumulated during the generation. If $p_i > p_j$, $i$'s strategy remains at $i$'s node in the next generation. However, if $p_i < p_j$, $j$'s strategy is copied onto $i$'s node for the next generation. In the event that $p_i = p_j$, a coin toss determines the prevailing strategy. As strategies are copied to the next generation each strategy component of every agent is subject to mutation with a probability $m = 0.10$. If selected for mutation, Gaussian noise with mean = 0 and standard deviation ± 0.01 is added to the component. Should mutation drive a component's value outside [0,1] the value is adjusted back to the closer boundary value.

*Introducing second-order punishment*
      In a second simulation, second-order punishment was introduced and agents were given the ability to punish observers who were too lenient on cheaters. In a CPD game with observer $k$, a new agent $l$ makes an assessment of whether $k$'s definition of a cheater is more lenient than $l$'s. It does this by determining whether $k$'s threshold for punishment $t_k$ is greater than its own $t_l$. If and only if $t_k > t_l$ then $l$ inflicts second-order punishment on $k$, and $l$ pays an amount $c_l$ to have $c_l M$ deducted from $k$'s net payoffs.

*Introducing retaliation*

The third simulation examined the effect of retaliation on cooperative outcomes. The base case simulation was modified so that an agent *i* automatically retaliated after being punished by paying an amount $s_i \in [0,1]$ to have its punisher sanctioned by the amount $s_i M$. Since $s_i$ could evolve to 0, agents might evolve so that they did not retaliate, even when punished. Three different rules were implemented for calculating how much a punished agent spent on retaliation. All methods of retaliation are arbitrary in the sense that their construction was intentionally limited to existing parameters of the model, but are nonetheless intuitive given the constraint of available variables. The three rules are:

1. $s_i$ equals the same amount the punished agent would have spent to punish a low contributor ($s_i = c_i$). This assumes that a single strategy component dictates how much an agent spends to punish others regardless of the reason.

2. $s_i$ is an independently evolving strategy component ($s_i$ is independent of $c_i$). This assumes that retaliation is a separate form of punishment and governed by its own strategy component.

3. $s_i$ equals the amount the agent contributes to the public good in the CPD ($s_i = x_i$). This reflects the idea that both punishment and public good contributions are non-selfish behaviors, and so may be governed by the same strategy component.

*Simulation variables and output*

The important parameter governing the mechanism of altruistic punishment is the ratio of costs incurred by the punishing party to those of the party being punished (Casari, 2005; Shutters, 2009). Defined above as the punishment multiplier *M*, this parameter is analogous to the strength or efficiency of punishment and, along with network type, is an independent variable in these simulations. The dependent variables of interest are the mean contribution and the mean payoff which evolve in a population after 10,000 generations. The mean contribution represents the population's level of cooperativeness while the mean payoff is a measure of the population's social welfare.

It is important to note that the magnitude of payoff values collected is somewhat arbitrary. A more meaningful measure is the magnitude of *change* in payoffs due to the various punishment and structural treatments. Thus, payoff results are presented in this study by a measure known as payoff efficiency, where 0% efficiency means that payoffs equal those expected in a population composed entirely of defectors without any form of punishment (6.0 in this case) and 100% means that all members of the population contribute their entire endowment to the public good and that no punishment of any kind takes place (for a mean payoff of 9.0 in this case). While it is not possible for a population to evolve higher than 100% payoff efficiency, it is possible for populations under punishment treatments to evolve negative payoff efficiencies. This is due to the additional costs incurred during acts of punishment, both by the punishee and the punisher.

For any given parameter set (Table 2), 100 replications were conducted at $M = 0.0$

and then at subsequent values of $M$ in increments of 0.5, up to $M = 10.0$. Because aggregate outcomes using retaliation still showed considerable variability when $M > 10.0$, simulations were run additionally from $M = 10.0$ to $M = 30.0$ in increments of 5.0.

**Table 2.** Simulation parameters and their values used in the continuous prisoner's dilemma.

| Parameter | Values |
|---|---|
| Population size ($N$) | 400 |
| Generations per simulation run | 10,000 |
| Games initiated by each agent in a generation | 3 |
| Value range for strategy components ($x$, $t$, $c$, $s$) | [0,1] |
| Probability of strategy component mutation ($m$) | 0.1 |
| Mean ±SD of Gaussian noise added during mutation | 0 ±0.01 |
| Punishment multiplier ($M$), [begin : end : step] | [0 : 10 : 0.5] |
| Public good multiplier ($r$) | 1.5 |

## Results and discussion

*Control case: effects of first-order punishment only*

In the first simulation agents played the CPD followed by a single round of punishment. Agents could pay $c$ to have a low contributor punished by an amount $cM$. This is the control case as neither second-order punishment nor retaliation was allowed. Consistent with previous studies (Gürerk et al., 2006; Shutters, 2012), cooperation evolved despite the fact that cooperators had no direct incentive to punish and could ignore cheaters without repercussions (Table 5). As $M$ increased in these simulations, cooperation evolved in all simulations with social structure (Figure 1, solid lines). For each regular network, at some threshold value of $M$ (Table 3) the population underwent a rapid transition from nearly full defection to nearly full cooperation. In simulations without social structure cooperation never evolved and mean contributions to the public good evolved to approximately 0.

Previous studies have suggested that altruistic punishment may only be sustained through group selection (Boyd et al., 2003). One may be puzzled then that this result exhibits sustained punishment without discrete groups. However, Wilson and Wilson (2007) assert that what is important for group selection is not that agents form discretely bounded groups, but that their social interactions are local compared to the entire population. This assertion is supported by the current results from simulations with simple (first-order) punishment only. Not only did punishment, and subsequently cooperation, emerge in networked populations where interactions are local, but the more localized, measured as lower average degree $d$, the more readily punishment proliferated (Table 3).

The value of $M$ at which populations transitioned to cooperation was particularly influenced by the mean degree of the network. This relates to a debate regarding the effect that network density has on the ability of a networked population to evolve cooperative behavior. Researchers have previously asserted that the more densely connected a population, the more likely that it will evolve cooperation (Marwell and Oliver, 1993; Opp and Gern, 1993), an assertion supported by Jun and Sethi's (2007) simulation experiment.

However, many recent studies suggest the opposite, showing that cooperation is inhibited in denser networks (Flache, 2002; Flache and Macy, 1996; Takács et al., 2008) and that increasing average degree requires increasing the relative benefit of cooperative acts before cooperation can emerge (Ohtsuki et al., 2006). Results from this study strongly support the latter view that denser networks *inhibit* the evolution of cooperation. Though full cooperation eventually evolved on all regular networks, the severity of punishment, *M*, required to evolve cooperative populations increased as the density of the network increased (Table 3). This finding is similar to that of Ifti et al. (2004) which showed that as neighborhood size increases beyond a critical threshold, cooperation collapses.

**Table 3.** Approximate value of *M* required for transition from defection to cooperation, without and with second-order punishment (2OP).

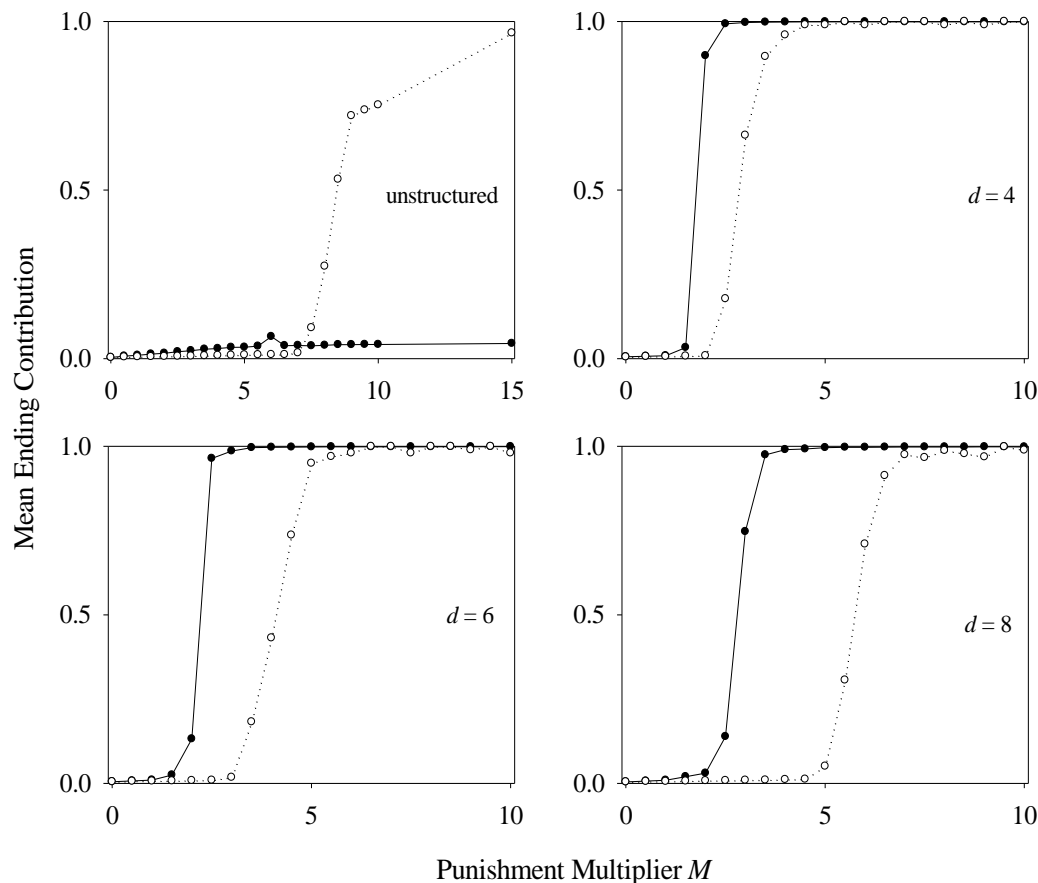| Network type | *d* | Approximate transition value of *M* | |
| --- | --- | --- | --- |
| | | Without 2OP | With 2OP |
| Complete | 399 | N/A[a] | 11.0 |
| Regular networks | 2 | 1.5 | 1.6 |
| | 4 | 1.8 | 2.8 |
| | 6 | 2.2 | 4.1 |
| | 8 | 2.8 | 5.7 |

[a] no transition occurred with increasing *M* even at values as high as $M = 5,000$.

*Effects of second-order punishment: Structured societies*

In the second set of simulations, agents could not only pay *c* to punish low contributors by an amount *cM*, they could also pay to punish those who had a higher tolerance for cheaters than themselves. Previous simulations have shown that when using a cultural group selection mechanism, second-order punishment may help to stabilize cooperative behavior in a population (Henrich and Boyd, 2001). However, results here show that instead of enhancing the cooperative effect of punishment, simulations using second-order punishment required higher values of *M* to induce cooperative behavior than simulations without social structure (Figure 1). In effect, punishment needed to be more severe to achieve cooperation than when there was no option for second-order punishment (Table 3).

One possible reason for this result is that in simulations with second-order punishment, agents that contributed fully to the public good could still suffer punishment for other reasons. Regardless of how cooperative they were, if they were lenient on cheaters, they might be the target of second-order punishment. Thus, many cooperative agents that might have helped move the population toward full cooperation could be injured through sanctions, making these punishers less fit and decreasing the overall effectiveness of punishment. This finding suggests that attempts to incite individuals to police each other through the threat of peer punishment may have unintended and adverse consequences.

**Figure 1.** Results of simulations with first-order punishment only and with both first- and second-order punishment.



*Note:* Results of simulations with first-order punishment only (solid lines) and with both first- and second-order punishment (dashed lines). Mean ending CPD contributions vs. *M* are presented for populations on a complete network and on three regular networks of varying degree. Each data point represents the mean ending contribution rate of 100 simulation runs.

It is important to understand that there are multiple ways to implement second-order punishment. In this study, an agent *l* bases its decision to inflict second-order punishment solely on an assessment of the *traits* of the observed first-order punisher *k*. Namely, *l* compares its own threshold for defining a cheater to the threshold of *k*. One alternative method of implementing second-order punishment is for *l* to observe the *behavior* of *k* in response to a third party *i*, where *i* is a participant in a CPD game. Once *k* determines whether or not to punish *i*, *l* then determines whether it would have taken the same action. If *k* reacted differently, then *l* inflicts second-order punishment on *k*. In other words, if *l* determines that *i* was a cheater and that *k* did not punish *i*, then *l* punishes *k*. Likewise, if *l* determines that *i* was a cooperator but was still punished by *k*, then *l* punishes *k* for being overly punitive. These last two cases may be implemented separately as well, leading to many alternative mechanisms for implementing second-order punishment. Therefore,

future research should seek to isolate the effects of different mechanisms of second-order punishment.

*Effects of second-order punishment: unstructured societies*

A surprising result was the ability of second-order punishment to induce cooperative outcomes in unstructured populations. Though simulations on complete networks evolved to full defection in every other treatment in this study, the addition of second-order punishment both increased cooperation and aggregate payoffs to relatively high levels (Table 4). With increasing *M*, levels of both public good contributions and payoffs using complete networks eventually surpassed those using regular networks (Figure 2).

This result suggests that at some point in a continuum of social structures, altruistic punishment alone becomes insufficient as a mechanism for upholding cooperation and second-order punishment emerges as a solution (see also Sigmund et al., 2010). If one considers this structural continuum as describing not simply the average degree of a society, but its overall size and complexity then a plausible example of the need for higher-order punishment can be viewed in the developmental dynamics of police agencies. As cities increase in population, and their policing agencies grow in size, the agencies typically add second-order punishment organizations (Wilson, 1963). Known variously as internal affairs, internal investigations, or similar designations, these organizations are responsible for policing the police. Evidence for this trend toward a need for higher order punishment may be further seen among the largest cities where citizen panels are frequently instituted to monitor the activities of internal affairs divisions. This emergence of third-order punishment may indicate that as societies continue along a continuum of societal size and complexity, increasingly higher order punishment levels are required to maintain cooperation.

**Table 4.** Mean ending payoff efficiency under different social structures, both without punishment and with different punishment treatments.

| Network type | *d* | No punishment | Punishment only[a] | Second-order punishment[a] | Retaliation (type 3)[a] |
|---|---|---|---|---|---|
| Complete | 399 | 0.003 | -0.007 | 0.883 | -0.043 |
| Regular networks | 2 | 0.006 | 0.541 | 0.157 | -0.090 |
| | 4 | 0.005 | 0.712 | 0.552 | -0.028 |
| | 6 | 0.005 | 0.736 | 0.663 | -0.060 |
| | 8 | 0.004 | 0.752 | 0.716 | -0.085 |

[a] mean ending payoff efficiency of 100 runs at $M = 15$

*Effects of retaliation*

In the third set of simulations, a punished agent was allowed to immediately retaliate against its punisher using one of three different rules (described above) to determine the amount *s* that the retaliating agent spent to impose a cost of *sM* on its punisher. Using retaliation rule 1 ($s_i = c_i$), cooperation did not evolve on any network. The ability to retaliate led to the collapse of cooperation that evolved when there was no

retaliation. Likewise, under rule 2 ($s_i$ is independent of $c_i$) full defection evolved on all social structures. In simulations using rule 3 ($s_i = x_i$) results were more complex. As with simple punishment, simulations with structured populations underwent a rapid transition from almost no contributions to some positive level of contributions with increasing $M$.

However, contributions did not transition to full cooperation as before but instead plateaued at a value between full cooperation and full defection, a value that varied by network density (Table 5). In addition, payoff efficiency initially rose with increasing $M$ but then fell to negative levels (Figure 3), meaning populations with the ability to retaliate fared worse than populations composed entirely of defectors and no punishment. Payoff efficiency decreased in the presence of retaliation even though some level of public good contribution was achieved. This result demonstrates the provisioning of public goods through what may be better described as coercion than cooperation.

Because humans often do retaliate after being punished (Molm, 1994; Nikiforakis, 2008), these results challenge the idea that cooperation may be the product of altruistic punishment in real world situations. If altruistic punishment actually has been an important mechanism in the evolution of cooperation, then it is likely that other mechanisms also existed to suppress or avoid retaliatory behavior. This may explain the frequency of institutional policies like those of the United States Department of Labor, which penalize or otherwise discourage retaliation against whistleblowers (USDL, 2009).
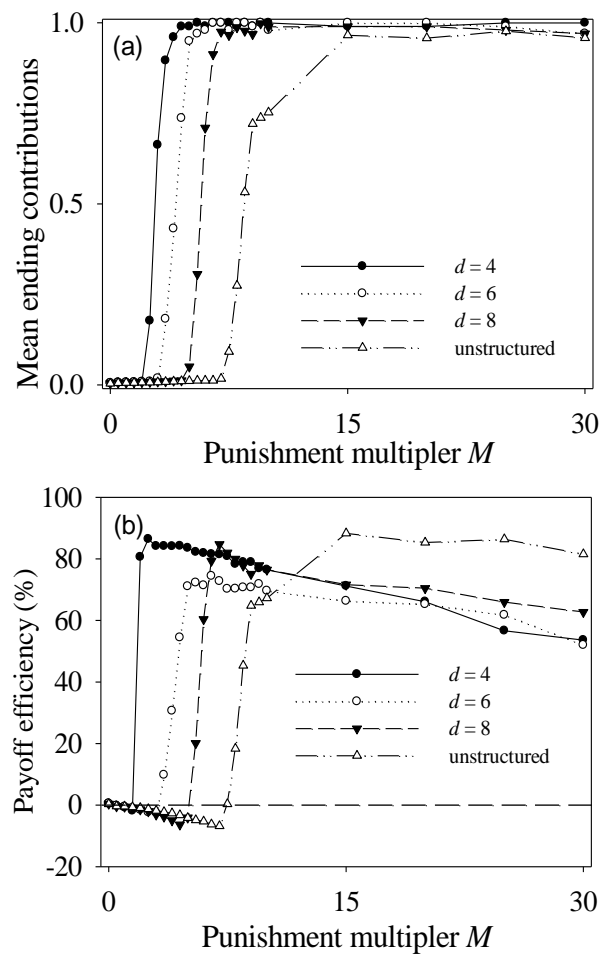
**Table 5.** Mean ending contributions ±SD under different social structures, without punishment, with one round of punishment, and with both a punishment and retaliation (rule 3) round.

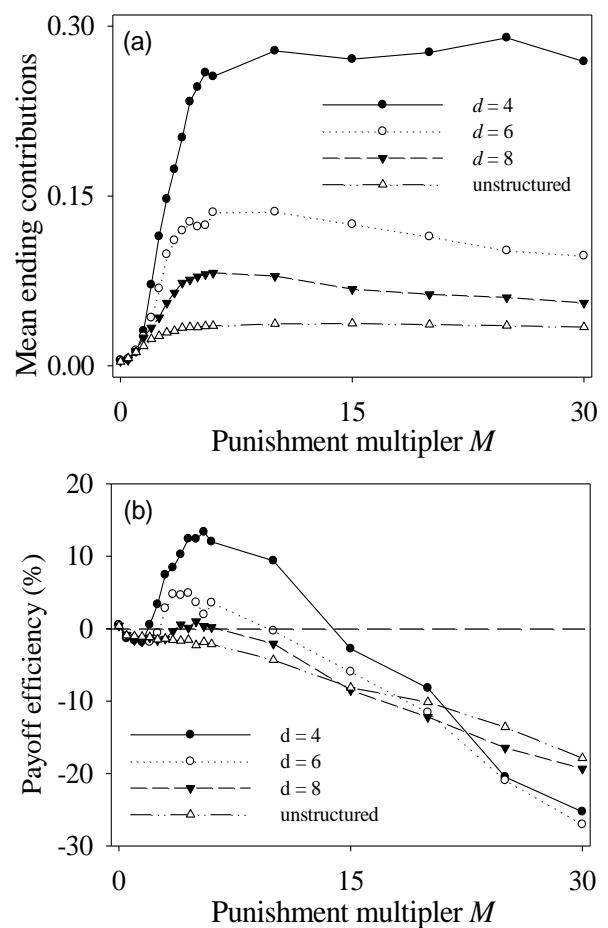| Network type | $d$ | No punishment | Punishment[a] | Retaliation[b] |
|---|---|---|---|---|
| Complete | 399 | 0.003 ±0.001 | 0.030 ±0.010 | 0.036 ±0.050 |
| Regular networks | 2 | 0.006 ±0.001 | 0.996 ±0.002 | 0.949 ±0.015 |
| | 4 | 0.005 ±0.001 | 0.998 ±0.001 | 0.277 ±0.064 |
| | 6 | 0.005 ±0.001 | 0.997 ±0.002 | 0.115 ±0.035 |
| | 8 | 0.004 ±0.001 | 0.990 ±0.017 | 0.065 ±0.018 |

[a] mean ending contribution of 100 runs at $M = 4$
[b] mean ending contribution of 100 runs at each $M = 10, 15, 20, 25, 30$ (500 total runs)

**Figure 2**. Effects of second-order punishment using four different networks.



**Figure 3**. Effects of retaliation using four different networks.



*Note:* (a) mean CPD contributions vs. *M* and (b) mean payoffs vs. *M* are presented.

*Note:* (a) mean CPD contributions vs. *M* and (b) mean payoffs vs. *M* are presented. Each point represents the population's mean contribution or payoff in the 10,000th generation of a single simulation run.

**Table 6.** Summary of effects of different punishment treatments on cooperative outcomes (public good contributions).

| Treatment type | Without social structure | With social structure |
| --- | --- | --- |
| Punishment of cheaters only | No cooperation emerged | Cooperation emerged at sufficiently high values of *M*, with the required value increasing as network density increased |
| Punishment of cheaters followed by second-order punishment | At sufficiently high values of punishment multiplier *M*, nearly full cooperation emerged | Cooperation emerged as with punishment of cheaters only, but required significantly higher values of *M* to emerge |
| Punishment of cheaters followed by retaliation | No cooperation emerged | Cooperation emerged as with punishment of cheaters only, but its magnitude was significantly lower, and decreased as network density increased |

*Further considerations of social structure*

In this study I have restricted structured populations to homogeneous regular networks to exclude confounding effects of variation among agents in degree, connectivity, edge effects, etc. However, regular networks bear little resemblance to the patterns of interactions among living things, particularly in humans, though they are arguably more representative of living systems than complete networks in which agents interact equally with all other members of a society. To briefly assess the effect of subsequent punishment behavior under more realistic social structures, supplemental simulations were run using small-world networks (Watts and Strogatz, 1998) and scale-free networks (Barabási, 2009; Tomassini, Pestelacci, and Luthi, 2007), both of which are common in complex physical and social systems (Barabási and Albert, 1999; Dorogtsev and Mendes, 2003).

Under small-world networks, results in all cases were qualitatively equivalent to results with regular networks presented in Table 6. However, results using scale-free networks present a challenge as neither second-order punishment nor retaliation appeared to have any effect on simulation outcomes. Both cases present ample opportunities for future research as they not only generate interesting results but are more applicable to the social structures under which social behavior likely evolved.

*Conclusion*

This study has built upon empirical studies that suggest altruistic punishment is a mechanism that leads to cooperation. Specifically, it examines two types of behavior that often occur in the presence of punishing behavior, retaliation and second-order punishment. Using computational social simulations, results show that retaliatory behavior almost always hinders the ability of punishment to induce cooperative behavior and that second-order punishment is most effective when populations are highly connected and/or well-mixed. These results concur qualitatively with observations from human social systems –

that retaliation is often suppressed and that second-order punishment frequently emerges when social systems grow beyond a certain threshold of size and complexity.

## References

Anderies, J. M. (2002). The transition from local to global dynamics: A proposed framework for agent-based thinking in social-ecological systems. In M. A. Janssen, (Ed.), *Complexity and ecosystem management* (pp. 13-34). Cheltenham, UK: Edward Elgar Publishing.

Axelrod, R., and Hamilton, W. D. (1981). The evolution of cooperation. *Science, 211*, 1390-1396.

Axelrod, R., Hammond, R. A., and Grafen, A. (2004). Altruism via kin-selection strategies that rely on arbitrary tags with which they coevolve. *Evolution, 58*, 1833-1838.

Barabási, A. L. (2009). Scale-free networks: A decade and beyond. *Science, 325*, 412-413.

Barabási, A. L., and Albert, R. (1999). Emergence of scaling in random networks. *Science, 286*, 509-512.

Barrett, S. (2003). *Environment and statecraft: The strategy of environmental treaty-making.* Oxford: Oxford University Press.

Boyd, R., Gintis, H., Bowles, S., and Richerson, P. J. (2003). The evolution of altruistic punishment. *Proceedings of the National Academy of Sciences of the United States of America, 100*, 3531-3535.

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

Boyd, R., and Richerson, P. J. (1989). The evolution of Indirect reciprocity. *Social Networks, 11*, 213-236.

Boyd, R., and Richerson, P. J. (2002). Group beneficial norms can spread rapidly in a structured population. *Journal of Theoretical Biology, 215*, 287-296.

Casari, M. (2005). On the design of peer punishment experiments. *Experimental Economics, 8*, 107-115.

Cederman, L. E. (2001). Agent-based modeling in political science. *The Political Methodologist, 10*, 16-22.

Chen, X. J., Fu, F., and Wang, L. (2007). Prisoner's dilemma on community networks. *Physica A: Statistical Mechanics and its Applications, 378*, 512-518.

Chwe, M. S. Y. (1999). Structure and strategy in collective action. *American Journal of Sociology, 105*, 128-156.

Cinyabuguma, M., Page, T., and Putterman, L. (2006). Can second-order punishment deter perverse punishment? *Experimental Economics, 9*, 265-279.

Clutton-Brock, T. H., and Parker, G. A. (1995). Punishment in Animal Societies. *Nature, 373*, 209-216.

Coleman, E. A., and Steed, B. C. (2009). Monitoring and sanctioning in the commons: An application to forestry. *Ecological Economics, 68*, 2106-2113.

Denant-Boemont, L., Masclet, D., and Noussair, C. N. (2007). Punishment, counterpunishment and sanction enforcement in a social dilemma experiment. *Economic Theory, 33*, 145-167.

Deng, K., and Chu, T. (2011). Adaptive evolution of cooperation through Darwinian dynamics in public goods games. *Plos One, 6*, e25496.

Dietz, T., Ostrom, E., and Stern, P. C. (2003). The struggle to govern the commons. *Science, 302*, 1907-1912.

Dorogtsev, S. N., and Mendes, J. F. (2003). *Evolution of networks: From biological nets to the Internet and WWW*. Oxford: Oxford University Press.

Dreber, A., Rand, D. G., Fudenberg, D., and Nowak, M. A. (2008). Winners don't punish. *Nature, 452*, 348-351.

Fehr, E., and Gächter, S. (2000). Cooperation and punishment in public goods experiments. *American Economic Review, 90*, 980-994.

Fehr, E., and Rockenbach, B. (2003). Detrimental effects of sanctions on human altruism. *Nature, 422*, 137-140.

Flache, A. (2002). The rational weakness of strong ties: Failure of group solidarity in a highly cohesive group of rational agents. *Journal of Mathematical Sociology, 26*, 189-216.

Flache, A., and Macy, M. W. (1996). The weakness of strong ties: Collective action failure in a highly cohesive group. *Journal of Mathematical Sociology, 21*, 3-28.

Fletcher, J. A., and Zwick, M. (2004). Strong altruism can evolve in randomly formed groups. *Journal of Theoretical Biology, 228*, 303-313.

Foster, K. R., and Ratnieks, F. L. W. (2001). Convergent evolution of worker policing by egg eating in the honeybee and common wasp. *Proceedings of the Royal Society of London. Series B: Biological Sciences, 268*, 169-174.

Frank, S. A. (1995). Mutual policing and repression of competition in the evolution of cooperative groups. *Nature, 377*, 520-522.

Gächter, S., Renner, E., and Sefton, M. (2008). The long-run benefits of punishment. *Science, 322*, 1510.

Goodnight, C. J. (2005). Multilevel selection: The evolution of cooperation in non-kin groups. *Population Ecology, 47*, 3-12.

Gould, R. V. (1993). Collective action and network structure. *American Sociological Review, 58*, 182-196.

Gürerk, Ö., Irlenbusch, B., and Rockenbach, B. (2006). The competitive advantage of sanctioning institutions. *Science, 312*, 108-111.

Hagen, E. H., and Hammerstein, P. (2006). Game theory and human evolution: A critique of some recent interpretations of experimental games. *Theoretical Population Biology, 69*, 339-348.

Hales, D. (2005). Change your tags fast! A necessary condition for cooperation? In P. Davidsson, B. Logan, and K. Takadama (Eds.), *Multi-agent and multi-agent-based simulation* (pp. 89-98). Berlin: Springer.

Hamilton, W. D. (1964). Genetical evolution of social behaviour 1 & 2. *Journal of Theoretical Biology, 7*, 1-52.

Harrison, N., and Singer, J. D. (2006). Complexity is more than systems theory. In N. Harrison (Ed.), *Complexity in world politics: Concepts and methods of a new paradigm* (pp. 25-41). Albany, NY: State University of New York Press.

Henrich, J., and Boyd, R. (2001). Why people punish defectors: Weak conformist transmission can stabilize costly enforcement of norms in cooperative dilemmas. *Journal of Theoretical Biology, 208*, 79-89.

Herrmann, B., Thöni, C., and Gächter, S. (2008). Antisocial punishment across societies. *Science, 319*, 1362-1367.

Hodgson, G. M. (2009). On the institutional foundations of law. *Journal of Economic Issues, 43*, 143-166.

Huang, Z. G., Wang, S. J., Xu, X. J., and Wang, Y. H. (2008). Promote cooperation by localised small-world communication. *Europhysics Letters, 81*, 28001.

Hui, P. M., Xu, C., and Zheng, D. F. (2007). Cooperation in evolutionary snowdrift game: Networking effects. *International Journal of Modern Physics B, 21*, 4035-4040.

Ifti, M., Killingback, T., and Doebelic, M. (2004). Effects of neighbourhood size and connectivity on the spatial continuous prisoner's dilemma. *Journal of Theoretical Biology, 231*, 97-106.

Jun, T., and Sethi, R. (2007). Neighborhood structure and the evolution of cooperation. *Journal of Evolutionary Economics, 17*, 623-646.

Kaul, I., and Mendoza, R. U. (2003). Advancing the concept of public goods. In I. Kaul, P. Conceição, K. L. Goulven, and R. U. Mendoza (Eds.), *Providing global public goods: Managing globalization* (pp. 78-111). New York: Oxford University Press.

Killingback, T., and Doebeli, M. (1996). Spatial evolutionary game theory: Hawks and doves revisited. *Proceedings of the Royal Society of London. Series B: Biological Sciences, 263*, 1135-1144.

Killingback, T., and Doebeli, M. (2002). The continuous prisoner's dilemma and the evolution of cooperation through reciprocal altruism with variable investment. *American Naturalist, 160*, 421-438.

Killingback, T., and Studer, E. (2001). Spatial ultimatum games, collaborations and the evolution of fairness. *Proceedings of the Royal Society of London. Series B: Biological Sciences, 268*, 1797-1801.

Leimar, O., and Hammerstein, P. (2001). Evolution of cooperation through indirect reciprocity. *Proceedings of the Royal Society of London. Series B: Biological Sciences, 268*, 745-753.

Marlowe, F. W., Berbesque, J. C., Barr, A., Barrett, C., Bolyanatz, A., Cardenas, J. C., . . ., and Tracer, D. (2008). More 'altruistic' punishment in larger societies. *Proceedings of the Royal Society of London. Series B: Biological Sciences, 275*, 587-590.

Marwell, G., and Oliver, P. E. (1993). *Critical mass in collective action.* Cambridge: Cambridge Publishing.

Maynard Smith, J., and Szathmáry, E. (1997). *The major transitions in evolution.* New York: Oxford University Press.

Michod, R. E. (1996). Cooperation and conflict in the evolution of individuality II: Conflict mediation. *Proceedings of the Royal Society of London. Series B: Biological Sciences, 263*, 813-822.

Molm, L. D. (1994). Is punishment effective? Coercive strategies in social-exchange. *Social Psychology Quarterly, 57*, 75-94.

Nikiforakis, N. (2008). Punishment and counter-punishment in public good games: Can we really govern ourselves? *Journal of Public Economics, 92*, 91-112.

Nikiforakis, N., and Engelmann, D. (2011). Altruistic punishment and the threat of feuds. *Journal of Economic Behavior & Organization, 78*, 319-332.

North, D. C. (2005). Understanding the process of economic change. Princeton, NJ: Princeton University Press.

Nowak, M. A., and May, R. M. (1992). Evolutionary games and spatial chaos. *Nature, 359*, 826-829.

Nowak, M. A., and Sigmund, K. (2005). Evolution of indirect reciprocity. *Nature, 437*, 1291-1298.

Ohtsuki, H., Hauert, C., Lieberman, E., and Nowak, M. A. (2006). A simple rule for the evolution of cooperation on graphs and social networks. *Nature, 441*, 502-505.

Opp, K. D., and Gern, C. (1993). Dissident groups, personal networks, and spontaneous cooperation  - The East-German revolution of 1989. *American Sociological Review, 58*, 659-680.

Ostrom, E., Walker, J., and Gardner, R. (1992). Covenants with and without a sword: Self-governance is possible. *American Political Science Review, 86*, 404-417.

Page, K. M., Nowak, M. A., and Sigmund, K. (2000). The spatial ultimatum game. *Proceedings of the Royal Society of London. Series B: Biological Sciences, 267*, 2177-2182.

Reeve, H. K., and Hölldobler, B. (2007). The emergence of a superorganism through intergroup competition. *Proceedings of the National Academy of Sciences of the United States of America, 104*, 9736-9740.

Richerson, P. J., and Boyd, R. (2005). *Not by genes alone: How culture transformed human evolution*. Chicago: University of Chicago Press.

Riolo, R. L., Cohen, M. D., and Axelrod, R. (2001). Evolution of cooperation without reciprocity. *Nature, 414*, 441-443.

Sandler, T. (1997). *Global challenges: An approach to environmental, political, and economic problems.* Cambridge, U.K.: Cambridge University Press.

Sandler, T. (1999). International cooperation and the international commons. *Duke Environmental Law & Policy Forum, 10*, 131-145.

Santos, F. C., Pacheco, J. M., and Lenaerts, T. (2006). Cooperation prevails when individuals adjust their social ties. *PLoS Computational Biology, 2*, 1284-1291.

Santos, F. C., Rodrigues, J. F., and Pacheco, J. M. (2006). Graph topology plays a determinant role in the evolution of cooperation. *Proceedings of the Royal Society of London. Series B: Biological Sciences, 273*, 51-55.

Sawyer, R. K. (2005). *Social emergence: Societies as complex systems.* Cambridge, UK:

Cambridge University Press.

Schofield, N. (1977). Dynamic games of collective action. *Public Choice, 30*, 77-105.

Sefton, M., Shupp, R., and Walker, J. M. (2007). The effect of rewards and sanctions in provision of public goods. *Economic Inquiry, 45*, 671-690.

Shutters, S. T. (2009). Strong reciprocity, social structure, and the evolution of fair allocations in a simulated ultimatum game. *Computational and Mathematical Organization Theory, 15*, 64-77.

Shutters, S. T. (2012). Punishment leads to cooperative behavior in structured societies. *Evolutionary Computation, 20*, 301-319.

Shutters, S. T., and Cutts, B. B. (2008). A simulation model of cultural consensus and persistent conflict. In V. S. Subrahamanian, and A. Kruglanski (Eds.), *Proceedings of the second international conference on computational cultural dynamics* (pp. 71-78). Menlo Park, CA: AAAI Press.

Sigmund, K., De Silva, H., Traulsen, A., and Hauert, C. (2010). Social learning promotes institutions for governing the commons. *Nature, 466*, 861-863.

Spector, L., and Klein, J. (2006). Genetic stability and territorial structure facilitate the evolution of tag-mediated altruism. *Artificial Life, 12*, 553-560.

Takács, K., Janky, B., and Flache, A. (2008). Collective action and network change. *Social Networks, 30*, 177-189.

Tomassini, M., Pestelacci, E., and Luthi, L. (2007). Social dilemmas and cooperation in complex networks. *International Journal of Modern Physics C, 18*, 1173-1185.

Travisano, M., and Velicer, G. J. (2004). Strategies of microbial cheater control. *Trends in Microbiology, 12*, 72-78.

Trivers, R. L. (1971). Evolution of reciprocal altruism. *Quarterly Review of Biology, 46*, 35-57.

USDL (2009). Whistleblower protections. *Summary of laws protecting whistleblowers from retaliation.* Retrieved March 6, 2009, from http://www.dol.gov/compliance/laws/comp-whistleblower.htm

Watts, D. J., and Strogatz, S. H. (1998). Collective dynamics of 'small-world' networks. *Nature, 393*, 440-442.

West, S. A., Griffin, A. S., and Gardner, A. (2007). Social semantics: Altruism, cooperation, mutualism, strong reciprocity and group selection. *Journal of Evolutionary Biology, 20*, 415-432.

Wilson, D. S., and Wilson, E. O. (2007). Rethinking the theoretical foundation of sociobiology. *Quarterly Review of Biology, 82*, 327-348.

Wilson, E. O., and Hölldobler, B. (2005). Eusociality: Origin and consequences. *Proceedings of the National Academy of Sciences of the United States of America, 102*, 13367-13371.

Wilson, J. Q. (1963). The police and their problems: A theory. *Public Policy, 12*, 189-216.