# Evolution of altruistic punishment in heterogeneous populations

Harmen de Weerd *, Rineke Verbrugge

Institute of Artificial Intelligence, Faculty of Mathematics and Natural Sciences, University of Groningen, PO Box 407, 9700 AK Groningen, The Netherlands

## ABSTRACT

Evolutionary models for altruistic behavior typically make the assumption of homogeneity: each individual has the same costs and benefits associated with cooperating with each other and punishing for selfish behavior. In this paper, we relax this assumption by separating the population into heterogeneous classes, such that individuals from different classes differ in their ability to punish for selfishness. We compare the effects of introducing heterogeneity this way across two population models, that each represents a different type of population: the infinite and well-mixed population describes the way workers of social insects such as ants are organized, while a spatially structured population is more related to the way social norms evolve and are maintained in a social network.

We find that heterogeneity in the effectiveness of punishment by itself has little to no effect on whether or not altruistic behavior will stabilize in a population. In contrast, heterogeneity in the cost that individuals pay to punish for selfish behavior allows altruistic behavior to be maintained more easily. Fewer punishers are needed to deter selfish behavior, and the individuals that punish will mostly belong to the class that pays a lower cost to do so. This effect is amplified when individuals that pay a lower cost for punishing inflict a higher punishment.

The two population models differ when individuals that pay a low cost for punishing also inflict a lower punishment. In this situation, altruistic behavior becomes harder to maintain in an infinite and well-mixed population. However, this effect does not occur when the population is spatially structured.

© 2011 Elsevier Ltd. All rights reserved.

## 1. Introduction

The question of how cooperation has evolved represents one of the more enduring puzzles in biology and social sciences, in which the role of many pieces is understood even if some pieces do not yet seem to fit together (Hamilton, 1964; Hardin, 1968; Axelrod and Hamilton, 1981; Sigmund, 2010; Gärdenfors, 2011). The paradox of cooperation is that although cooperation adds to the common good of a group of individuals, contributing to the common good generally bears a higher cost than the individual returns (Hardin, 1968; West et al., 2007). Individuals that enjoy the cooperation of others without being cooperative themselves are therefore at an evolutionary advantage. At first sight, a group of individuals thus seems to be destined to never cooperate, even if the combined benefit of every single individual cooperating outweighs the cost of contributing.

Even though cooperation seems to be destined to fail in theory, many social animals engage in cooperative action, ranging over a wide variety of activities (Wilkinson, 1984; Mulder and Langmore, 1993; Dugatkin, 1997; Crespi, 2001). To explain why cooperation stabilizes in many animal societies, a number of mechanisms have been proposed, varying in the assumptions they make on individual cognitive abilities (see among others Nowak, 2006; Gärdenfors, 2011). One of the mechanisms that may stabilize cooperation is punishment (Sigmund et al., 2001; Boyd et al., 2003; Brandt et al., 2003; Fowler, 2005). Punishment can provide the necessary incentive to stabilize cooperation in animal (Clutton-Brock and Parker, 1995; Monnin and Ratnieks, 2001) as well as in human societies (Ostrom, 2000; Fehr and Gächter, 2002). Experiments have shown that human subjects have a high willingness to sacrifice in order to punish selfish behavior, even when punishment is understood to yield no future benefits (Güth et al., 1982; Camerer and Thaler, 1995; Bolton and Zwick, 1995; Henrich et al., 2001; Fehr and Gächter, 2000).

An $N$-person extension of the prisoner's dilemma, known as the public goods game (Kagel et al., 1995; Fehr and Gächter, 2002), has been investigated in simulations of the evolution of cooperation. In the public goods game, the game is played by $N > 2$ individuals, each of which receives an initial capital $C$. They may choose to keep that capital to themselves, or invest any part of it in a common pool. Once every player has decided how much to invest, the capital in the common pool is doubled, and divided equally among the players, irrespective of their investment. If every player invests their entire

* Corresponding author. Tel.: +31 50 363 4114; fax: +31 50 363 6687.
E-mail addresses: hdeweerd@ai.rug.nl (H. de Weerd),
rineke@ai.rug.nl (R. Verbrugge).

capital, each will end up with $2C$ and therefore double their initial capital. However, each individual is faced with the temptation of exploiting the common pool. Since every individual investment is divided equally among all $N > 2$ individuals, the return on the individual investment is negative. The game-theoretical dominant strategy would therefore be to invest nothing. But if none of the players invests, each will end up with their initial capital $C$, which is half the capital they would have gained if everyone had invested. In experiments with volunteers with actual economic incentives, human players do tend to invest a reasonable sum. Typically, in the first round, participants choose to invest at least half their capital. When the game is repeated over several rounds, the amount invested quickly declines until nobody invests anything, unless there is an opportunity to punish individuals for low investments (Fehr and Gächter, 2002) or opt out of playing the public goods game (Orbell and Dawes , 1993; Semmann et al., 2003).

The models that have been proposed so far to explain why cooperation and punishment persist and would even be able to invade in a population of selfish individuals, commonly make the assumption of homogeneity. In a homogeneous society, individuals can use different strategies, but the payoffs of an encounter between two individuals depend only on the strategy the individuals adopt. Individuals have the same cost of punishing, and the same benefit of their partner cooperating. In this article, we propose to relax this assumption of homogeneity by allowing for populations that consist of two or more different sub-classes. By allowing the costs and benefits of cooperation to vary across sub-classes, individuals from different classes may have different opportunities. Empirical research shows that differences in marginal benefit from contributions to a public good changes the willingness to contribute and punish (Fisher et al., 1995; Reuben and Riedl, 2009). Subjects that enjoy a higher benefit not only tend to contribute more to the public good, but are also expected to do so, and are punished more severely by other players if they contribute less than their fair share.

In this research, we determine the effects of a heterogeneous population of individuals on the evolution of altruistic punishment, and the resulting structure of the population in a simulated environment. We adjust the model of the public goods game with voluntary participation introduced by Hauert et al. (2002, 2002) and further extended to include altruistic punishment (Fowler, 2005; Brandt et al., 2006) to allow for heterogeneous classes of individuals. Specifically, we investigate the effect of individual differences in the cost for punishing a co-player as well as the cost of being punished by another individual. We compare these effects across two different population models. In our first model, discussed in Section 2, the public goods game is played in an infinite size and well-mixed population, where individuals are assumed never to encounter each other more than once in the same setting. Section 3 describes the second model, which imposes a spatial structure on the population in the form of a lattice, such that individuals only interact with a small selection of close neighbors. For both population models, we present a model for a population that is divided into $M$ classes of individuals, and show the results of an implementation of the model for the case of $M=2$ classes. The individuals we simulate share the knowledge that the population is heterogeneous, but not how this affects the rewards. Simulation results are presented separately for each model, while Section 4 summarizes these results and provides directions for further research.

## 2. Infinite population model

To determine the effect of heterogeneity of individuals on the evolution of altruistic punishment, we have constructed two model variations of the public goods game. In this section, we will discuss a model based on the assumption of an infinite sized, well-mixed population of individuals. This model can be used to represent any sufficiently large population in which individuals are very unlikely to encounter the same co-player twice in the setting of a public goods game over the course of their lifetime. The infinite population model may therefore describe the public goods game in a large colony of social insects, such as ants, bees or wasps. In these societies, workers generally exhibit altruistic behavior by sacrificing most or all of their direct reproduction to help rear the offspring of the queen (Oster and Wilson, 1979). Interestingly, in some species of social insects, infertile workers can still lay haploid eggs destined to be males (Wenseleers et al., 2005). There is an evolutionary incentive to do so when the queen is mated to more than two males, in the sense that workers are more related to their own sons than to sons of their queen mother and sons of their sister workers (Trivers and Hare, 1976; Wenseleers et al., 2004). The reward for such behavior is therefore an increase in their inclusive fitness, that is the probability of their genes surviving. Natural selection would therefore favor the social insect that lays its own eggs. However, workers lay eggs at the expense of performing their duties to the colony. Punishment for this selfish behavior takes the form of queen and worker policing (Monnin and Ratnieks, 2001). Through this mechanism, worker-laid eggs are destroyed, effectively removing all benefits from the selfish behavior.

Even though workers, queens and males in a colony of social insects represent morphologically different castes that perform different tasks, the homogeneous infinite population model can be used to model the interactions between the workers of large colonies. However, some colonies of social insects exhibit a further subdivision of the worker caste (Oster and Wilson, 1979) up to a point where a heterogeneous infinite population model would fit the situation better. For example, leaf-cutting ant workers exhibit a 200-fold variation in body mass (Wilson, 1980), while in weaver ants of the genus *Oecophylla*, workers show a clear bimodal size distribution, with almost no overlap in size between minor and major workers (Hölldobler and Wilson, 1990). In cases like these, morphologically different workers typically perform different tasks depending on their physical traits. In general, the minor workers stay in the nest to tend to the queen and her brood, while major workers perform the more dangerous tasks of foraging and defending the colony (Hölldobler and Wilson, 1990).

In the remainder of this section, we will discuss how heterogeneity between individuals affects the evolution of altruistic punishment in the infinite population model. As a starting point, we use the model introduced by Brandt et al. (2006), which already allows for voluntary participation. This model is extended in the present work by dividing the population into $M$ heterogeneous classes of individuals. For the simulation results, we restrict ourselves to the case $M=2$.

### 2.1. Infinite population model: methods

In the infinite population model, we follow Brandt et al. (2006). Their model is an extension of the basic public goods model, in which players may choose not to share in the public good and instead receive a fixed payoff. We further extend their model to allow for heterogeneous groups within the population. In our case, the population is assumed to consist of $M$ classes of individuals, which occur at a fixed ratio within the population. That is, although evolutionary dynamics affect the frequencies at which the different strategies occur within each class, this has no effect on the relative frequency of the different classes within the population. In effect, this means there is no genetic basis that determines the individual membership to a class. Each class of individuals occurs at a constant frequency $0 < f_i < 1$ ($1 \leq i \leq M$), such that $\sum_i f_i = 1$.

Each class of individuals $i$ is further divided by the strategy they adopt: the loners $x_{i,L}$, altruistic non-punishers or cooperators $x_{i,AN}$, selfish non-punishers or defectors $x_{i,SN}$ and altruistic punishers $x_{i,AP}$, where $x_{i,L}$, $x_{i,AN}$, $x_{i,SN}$, $x_{i,AP}$ refer to the fraction of the population adopting their respective strategy such that

$$x_{i,AN} + x_{i,AP} + x_{i,SN} + x_{i,L} = f_i \quad \text{for all } 1 \leq i \leq M.$$

For convenience, the notation $x_s$ is used to denote the fraction of the entire population that adopts strategy $s$. That is,

$$x_s := \sum_{i=1}^{M} x_{i,s} \quad \text{for all strategies } s \in \{AN, AP, SN, L\}.$$

Note that in this setting individuals only play pure strategies. We assume that the public goods game is not played by the entire population simultaneously. Instead, the game is played by a random sample of $N$ individuals. The expected payoffs of each of the strategies are calculated accordingly.

When a random sample of size $N$ is drawn, the public goods game is played by all the individuals in this group except for the loners. Loners refuse to play the game and instead of sharing in the public goods, they receive a fixed payoff $\sigma$. They have no share in the public good, but they also do not contribute to it, and are not punished for failing to contribute.

Among the individuals that decide to play the game, altruistic individuals choose to invest an amount $c$ in the public goods. The total amount contributed in the public goods by all of the $N$ individuals is multiplied by a factor $r > 1$ before it is distributed among all individuals playing the game, whether they are altruistic or selfish, but excluding the loners. That is, in a group of $n_A := N(x_{AP} + x_{AN})$ altruistic individuals and $n_L := Nx_L$ loners, the public goods yield the non-loners a benefit of $rc \cdot n_A/(N - n_L)$ at a cost $c$ to each of the altruistic individuals. However, there is an exception to this rule. When the group consists of $N - 1$ loners, the only individual willing to participate in the public goods game is forced to be a loner as well. Furthermore, it is assumed that $(r-1)c > \sigma > 0$, such that a loner receives a better payoff than the members of a group of selfish individuals that receive 0 payoff, but worse than the members of a group of altruistic individuals, where each individual receives $(r-1)c$.

After all contributions have been made and the public good is shared, punishing individuals punish the selfish individuals. Selfish individuals in this setting are individuals that choose to participate in the public goods game, but do not contribute to the public good. The punishment they receive for this behavior depends on the class of the altruistic punisher. Altruistic punishers of class $i$ inflict a cost $\beta_i > 0$ to each selfish individual at a personal cost of $\gamma_i > 0$. Following Brandt et al. (2006), altruistic punishers furthermore punish individuals that fail to punish selfish participants. In this context, altruistic non-punishers are also termed second-order free-riders, since they do contribute to the public good, but do not contribute to the punishment system. Altruistic punishers in class $i$ punish second-order free-riders for a fraction $0 \leq \alpha \leq 1$ of the punishment they inflict to selfish individuals. That is, at a cost of $\alpha\gamma_i$ to themselves, they lower the payoff of altruistic non-punishers by $\alpha\beta_i$. However, when there are no selfish individuals in the group, none of the participants of the public goods game will punish, and altruistic non-punishers can therefore not be detected. Second-order free-riding is therefore only punished if there are at least one altruistic punisher, at least one altruistic non-punisher, and at least one selfish individual present in the group.

Evolution in well-mixed, infinite populations is traditionally studied using replicator dynamics (Taylor and Jonker, 1978; Nowak and Sigmund, 2004). We follow the method outlined in Brandt et al. (2006) to determine the expected payoffs $P_{i,s}$ for individuals of class $i$ that adopt strategy $s$. This results in the following expected payoffs:

$$P_{i,L} = \sigma,$$

$$P_{i,AN} = \sigma x_L^{N-1} + rc(x_{AP} + x_{AN})B(x_L) - cF(x_L) - \alpha G(x_{SN}) \sum_{j=1}^{M} \beta_j(N-1)x_{j,AP},$$

$$P_{i,SN} = \sigma x_L^{N-1} + rc(x_{AP} + x_{AN})B(x_L) - \sum_{j=1}^{M} \beta_j(N-1)x_{j,AP},$$

$$P_{i,AP} = \sigma x_L^{N-1} + rc(x_{AP} + x_{AN})B(x_L) - cF(x_L) - \gamma_i(N-1)(x_{SN} + \alpha x_{AN}G(x_{SN})),$$

where the auxiliary functions $B$, $F$ and $G$ are defined analogously to Brandt et al. (2006):

$$B(x_L) = \frac{1}{1-x_L}\left(1 - \frac{1-x_L^N}{N(1-x_L)}\right)$$

$$F(x_L) = 1 + x_L^{N-1}(r-1) - \frac{r}{N}\frac{1-x_L^N}{1-x_L}$$

$$G(x_{SN}) = 1 - (1-x_{SN})^{N-2}.$$

Note that in our heterogeneous setting, the class of an individual only affects the cost for punishing ($\gamma_i$) as well as the effectiveness of punishment ($\beta_i$). The costs and benefits of participating in the public goods game, as well as the loner payoff, are the same for all individuals in the population.
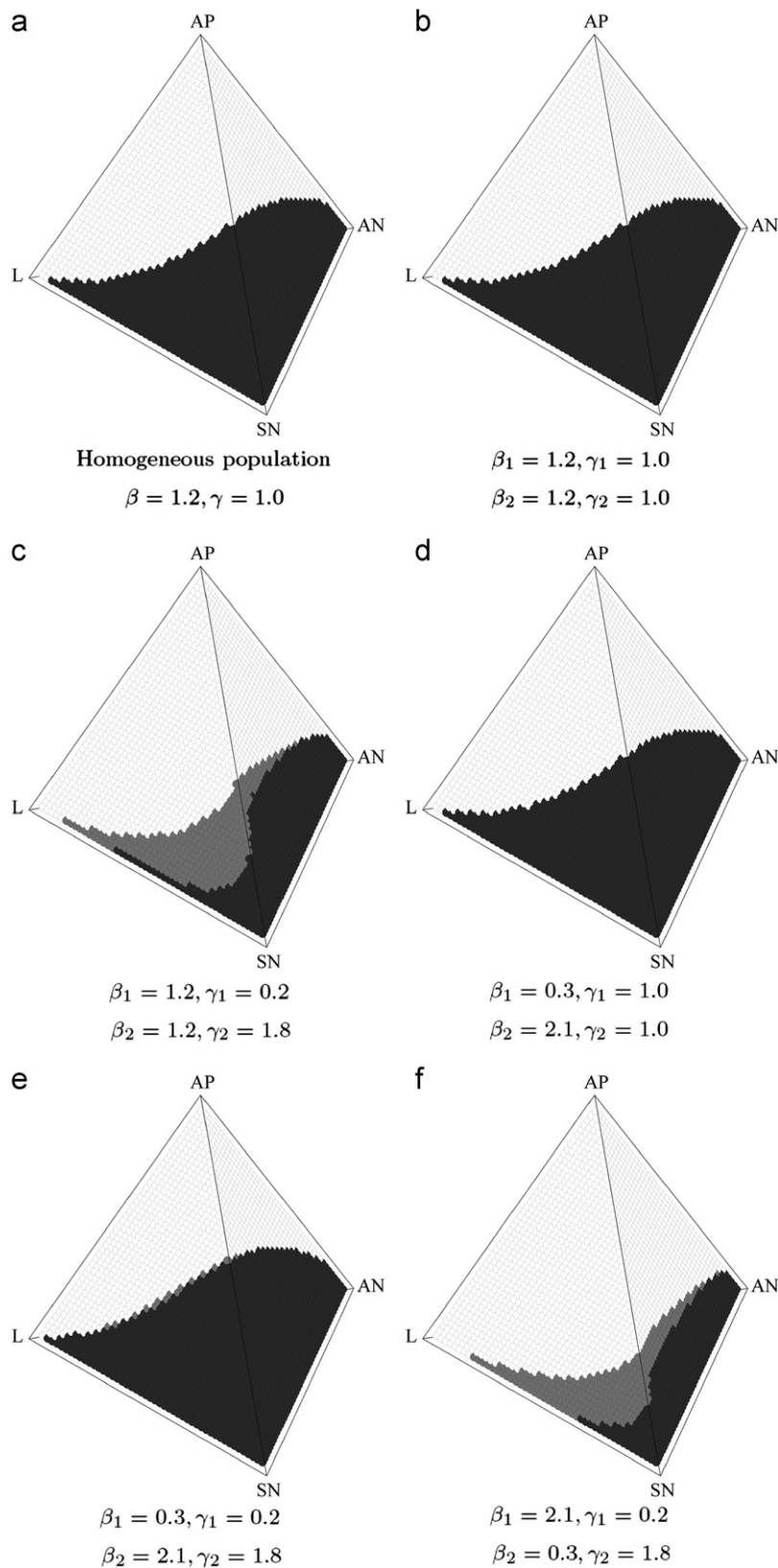
## 2.2. Infinite population model: results

To determine the effects of heterogeneity of individuals within a population on the behavior of a well-mixed infinite population, the model outlined in Section 2.1 has been implemented in Java. Following Brandt et al. (2006), we use the parameter setting $c = \sigma = \gamma = 1.0, \beta = 1.2, \alpha = 0.1, r = 3.0, N = 5$.

Since the proportions of the population that have adopted the strategy $AN$, $SN$, $L$ and $AP$ sum up to one, the configuration of strategies in a single homogeneous population ($M=1$) can be represented as a point in the simplex $S_4$; the convex hull of the pure strategies $AN$, $SN$, $L$ and $AP$. Each point $p$ within the simplex represents a configuration for which the relative frequency of each strategy is proportional to the distance of $p$ to the corresponding corner. Similar to Brandt et al. (2006), we find that any point in the interior of the simplex is drawn to either one of two *attractors*.[1] That is, given enough time, any population will eventually settle into one of the two possible situations. Which situation the population will end up in is only determined by the initial proportions of the different strategies. Since the interactions between the proportions of the four strategies are no longer mathematically tractable, Fig. B.1a shows the results of numerical simulations. Each dot indicates an initial state of the population in the interior of the simplex. That is, for each dot in Fig. B.1a, each of the proportions $x_{SN}, x_{AN}, x_{AP}$ and $x_L$ is initially non-zero. The brightness of the dot indicates the final destiny of the population. The first possibility, indicated by dark grey dots in Fig. B.1a, is the set of periodic orbits in the plane $x_{AP} = 0$. In this case, all punishers are eliminated from the population, and the population settles in a periodic orbit in the AN-SN-L plane. The second attractor is the AP-AN edge, in which the population only consists of altruistic individuals. The initial states that end up in this situation are indicated by white dots in Fig. B.1a.

To determine the effects of heterogeneity, we introduced the heterogeneous classes of individuals as described in Section 2.1.

---

[1] See Brandt et al. (2006) for a complete discussion of the results in a homogeneous population.

**a**

AP

AN

L

SN

Homogeneous population
$\beta = 1.2, \gamma = 1.0$

**b**

AP

AN

L

SN

$\beta_1 = 1.2, \gamma_1 = 1.0$
$\beta_2 = 1.2, \gamma_2 = 1.0$

**c**

AP

AN

L

SN

$\beta_1 = 1.2, \gamma_1 = 0.2$
$\beta_2 = 1.2, \gamma_2 = 1.8$

**d**

AP

AN

L

SN

$\beta_1 = 0.3, \gamma_1 = 1.0$
$\beta_2 = 2.1, \gamma_2 = 1.0$

**e**

AP

AN

L

SN

$\beta_1 = 0.3, \gamma_1 = 0.2$
$\beta_2 = 2.1, \gamma_2 = 1.8$

**f**

AP

AN

L

SN

$\beta_1 = 2.1, \gamma_1 = 0.2$
$\beta_2 = 0.3, \gamma_2 = 1.8$

**Fig. B.1.** Interior of the simplex $S_4$ for six different settings: (a) single homogeneous population, (b) homogeneous classes, (c) low cost, (d) low returns, (e) low returns and low cost, and (f) high returns and low cost. To improve comparability across situations, the average values $\overline{\beta}$ and $\overline{\gamma}$ over the two classes in situations (b)–(f) are fixed. Dark grey dots indicate initial configurations of the population for which individuals of both classes are drawn to the plane $x_{AP} = 0$, while configurations drawn to the AP-AN edge are marked by white dots. When classes are drawn to different attractors, the corresponding point in the figure is light grey.

That is, the population is subdivided into $M > 1$ classes such that the return on punishment $\beta_i$ and the cost of punishment $\gamma_i$ are homogeneous within a class, but heterogeneous across classes.

For our simulation, the number of classes was limited to $M=2$. The average values $\overline{\beta}$ and $\overline{\gamma}$ were fixed such that the average effectiveness of punishment and average cost of punishing are the

same for the homogeneous and the heterogeneous populations in the initial configuration. Fig. B.1b–B.1f shows simulation results[2] in the case each class represents half the population, and the initial configuration of strategies is constant across classes. That is, for each strategy $s \in \{SN, AN, AP, L\}$, $x_{1,s} = x_{2,s}$ at the beginning of the simulation. This way, the results can be presented in a simplex similar to the homogeneous case. Like before, initial configurations for which the entire population is drawn to the AP-AN edge are represented by a white dot in Fig. B.1. Similarly, a dark grey dot indicates that an initial configuration is drawn to the AN-SN-L plane. A new situation arises when for some initial configuration, different classes of individuals within the population are drawn to different attractors. In Fig. B.1, such initial configurations are indicated by a light grey dot.

For completeness, Fig. B.1b shows the situation of a population consisting of two homogeneous classes. In this setting, the population is subdivided into two classes, but the individuals of the different classes are completely homogeneous. Conceptually, this setup should produce the same results as a population consisting of only one class of homogeneous individuals, although numerical simulation could introduce some differences. Instead of a single homogeneous population, the population is divided into two classes consisting of indistinguishable individuals that can only learn from other individuals in the same class. However, since the initial state of each class is the same, both classes react exactly the same.

To determine the effects of heterogeneity, we first consider the situation in which the cost of punishing $\gamma_i$ is taken to be heterogeneous across classes, but the effectiveness of punishment $\beta_i$ is constant across classes. In this case, payoffs are dependent on the individual's class, since altruistic punishers of class 1 will pay a lower cost for punishing selfishness ($\gamma_1 = 0.2$) while individuals adopting the same strategy in class 2 will pay a higher cost ($\gamma_2 = 1.8$). Although each class of individuals still has the same two attractors as in the case of the homogeneous population, different classes may be drawn to different attractors. Individuals of class 1 may find it lucrative to punish in the sense that punishing reduces the fitness of selfish individuals more than it reduces the fitness of the punishers themselves, while for class 2 the penalty incurred for selfishness may be too low when offset against the costs of inflicting such a penalty.

Fig. B.1c shows the effect of heterogeneous cost of punishment on the final state of the population. The light grey dots indicate initial configurations for which class 1, having the lower cost of punishment $\gamma_1 = 0.2$, is drawn to the AP-AN edge, while the higher cost of punishment $\gamma_2 = 1.8$ causes all altruistic punishers to disappear from class 2. In the simulation, each of these situations resulted in a population of only altruists, under the influence of the altruistic punishers of class 1. In effect, in these cases the burden of punishing selfishness is carried by the individuals best suited for the task in the sense that they are the more efficient punishers.

One effect that appears in Fig. B.1c is that under heterogeneous classes, an initial configuration will end up in an end state without punishers more readily when the proportions of selfish non-punishers and altruistic non-punishers are balanced. This is due to the fact that altruistic non-punishers will contribute to the common good, which helps selfish non-punishers and altruistic punishers equally in terms of fitness. In effect, altruistic non-punishers compensate selfish non-punishers for the punishment they receive.

When instead of the cost of punishing $\gamma_i$ we consider the situation in which the returns on punishment $\beta_i$ vary across classes, the picture changes. Note that changes in $\beta_i$ only affect the payoff of

selfish individuals and, more importantly, uniformly so for all classes. In this case the payoff of an individual only depends on the strategy it adopts, and is independent of its class. Therefore, the results for this heterogeneous population model are similar to the homogeneous population model. Whether an individual will be drawn to a final state with only altruistic individuals (the AP-AN edge) or tend to a solution without punishers (the AN-SN-L plane) is also independent of class. There is no configuration which leads to different classes being drawn to different attractors. Moreover, if the initial proportions of strategies are the same for each class, the population as a whole will react exactly the same as a homogeneous class of individuals with $\beta = \sum f_i \beta_i$, where $f_i$ denotes the relative frequency of individuals belonging to class $i$. This result is shown in Fig. B.1d. Even though the two classes in this simulation differ in the effectiveness of their punishment ($\beta_1 = 0.3$ and $\beta_2 = 2.1$), since their (weighted) average effectiveness is the same as the case with homogeneous classes shown in Fig. B.1b, there are no real differences between the figures.

Even though heterogeneity in the returns on punishment $\beta_i$ between classes does not influence the behavior of the population by itself, it does change the behavior of the population when the cost of punishing $\gamma_i$ is also heterogeneous. Fig. B.1e shows that the interaction between the two types of heterogeneity can deter altruistic punishment. In this situation, class 1 inflicts a low punishment at a low cost ($\gamma_1 = 0.2$ and $\beta_1 = 0.3$), while class 2 inflicts high punishment at a high cost ($\gamma_1 = 1.8$ and $\beta_1 = 2.1$). When classes differ in the level of punishment they inflict this way, Fig. B.1e shows that the proportion of final states that include punishers falls below the baseline performance of the homogeneous population. In an infinite and well-mixed population, punishing at different levels may hinder the evolution of altruistic punishment.

Fig. B.1f shows that the interaction between heterogeneity in cost for and returns on punishment can work both ways. When class 1 has a lower than average cost of punishment $\gamma_1 = 0.2$ and also a higher than average return on punishment $\beta_1 = 2.1$, this increases the proportion of initial configurations for which the final state includes punishers. This effect is stronger than the separate effect of heterogeneity in the costs for punishing, despite the fact that class 2 is virtually ineffective at punishing for selfishness, with altruistic punishers paying a cost of $\gamma_2 = 1.8$ to inflict a punishment of $\beta_2 = 0.3$.

## 2.3. Infinite population model: discussion

We have shown that for the infinite population model, a population can take advantage of heterogeneity in the ability to punish for selfish behavior by specialization. When the cost for punishing is differentiated, punishing co-players for selfish behavior can be feasible for some class of individuals, while it may be too expensive for another class. Heterogeneity in the cost for punishing makes it easier for altruistic behavior to evolve; a lower proportion of altruistic punishers is needed to ensure that a population will end up in a state without selfish individuals. In situations in which the difference in cost is high, or when there are many individuals exhibiting selfish behavior, this may lead to a clear specialization. Altruistic punishers disappear from the class with the highest cost for punishing, which leaves the class with the lowest cost for punishing with the responsibility of enforcing altruistic behavior through punishment. The model therefore predicts that when individual differences in the costs of punishing are sufficiently high, these differences will cause a division of labor in a well-mixed population.

In contrast, the infinite and well-mixed population proved insensitive to heterogeneity in the returns on punishment. When controlled for the average value, variations in the returns on punishment across classes of individuals do not change the way altruistic behavior and punishment evolve. However, variations in

---

[2] Results for the entire interior of the simplex are available online at http://www.ai.rug.nl/SocialCognition/?p=83

the returns on punishment do interact with heterogeneity of punishing cost. If some class of individuals can inflict high punishment, doing so at a low cost further increases the chances for altruistic behavior to stabilize. If, on the other hand, high punishment can only be achieved by a high cost, the combined effect will make it harder for the population to stabilize altruistic behavior than in the case of a homogeneous population. The model therefore predicts that large differences in the level of punishment are rare in well-mixed populations, and that when individuals that pay a lower cost for punishing for selfishness also inflict a higher punishment, these individuals are likely to become solely responsible for punishing.

## 3. Spatial model

In Section 2, we derived a model for playing the public goods game under the assumption of an infinite size, well-mixed population of individuals. From simulations with repeated pairwise interactions between individuals, such as the prisoner's dilemma or the ultimatum game, it is known that including spatial structure to the model of the population can have strong effects on the evolution of cooperation and punishment (Nowak and May, 1993; Lindgren and Nordahl, 1994; Killingback and Doebeli, 1996). The spatial structure allows altruistic individuals to persist through positive assortment, locally avoiding exploitation from selfish individuals by clustering together. In this section, we will derive a model for playing the public goods game in a spatially structured environment, and extend it to allow for heterogeneous classes of individuals.

Unlike in the infinite population model, in which individuals indiscriminately interact with every other individual they encounter, each individual in the spatial model only interacts with a specific group of other individuals. In this sense, the spatial model can be used to represent models such as kin or group selection. In these models, individuals preferentially interact with specific individuals, either because of relatedness, or more pragmatic reasons such as spatial distance (Hamilton, 1964; Griffin and West, 2002; West et al., 2002; Fletcher and Zwick, 2007; Lehmann et al., 2007; Nowak et al., 2010).

In the remainder of this section, we discuss a spatial model of playing the public goods game (inspired by Hauert et al., 2002, 2002), and extend it in the present work to allow for the population to be subdivided into $M$ heterogeneous classes. To determine the effects of heterogeneity on the evolution of altruistic punishment, we present simulation results of this model, in which the number of classes is restricted to $M=2$.

### 3.1. Spatial model: methods

In order to keep the results of the spatially structured world comparable to the results of the infinite well-mixed population model of Section 2, we chose to organize the individuals on a square lattice similar to Hauert et al. (2002). In the lattice setting, interactions between individuals are limited to include only those within a certain spatial neighborhood. To prevent edge effects, periodic boundaries are assumed. That is, the lattice represents a torus, such that the left edge of the lattice is first connected to the right edge, and the top edge is then connected to the bottom.

The size of the interaction neighborhood affects the eventual state of the population. Ifti et al. (2004) show that in the Continuous Prisoner's Dilemma, in which cooperation is measured as an amount invested in cooperation rather than a binary choice, smaller neighborhoods tend to favor cooperation, while cooperation becomes unsustainable when the interactions are possible over larger distances. Ifti et al. (2004) also report that

when the learning and interaction neighborhood differ in size, the final state of any population playing the Continuous Prisoner's Dilemma game is zero cooperation. Even though our models differ from the one used in Ifti et al. (2004), we choose to keep the interaction and learning neighborhood of equal size. That is, every individual will compare its fitness only with individuals that it played the public goods game with during that round. As with the infinite and well-mixed population of Section 2, individuals only compare their fitness to the fitness of individuals with which they share a class.

In the spatial model we investigated, a public goods game is played by an individual and its four direct neighbors.[3] Thus, all public goods games played have a maximum of five participants if none of them chooses to be a loner. Also, since each of the individuals "hosts" one game, each individual plays the public goods game a total of five times. The game is divided into discrete rounds. After each round, all individuals simultaneously update their strategy by adopting the strategy that yielded the highest fitness among those individuals they interacted with during that round.

In the spatial model, the payoffs of the public goods game are actual payoffs rather than expected payoffs. We define $n_{i,s}$ as the number of individuals in the group that are of class $i$ and have adopted strategy $s$, and $n_s$ as the total number of individuals that have adopted strategy $s$ such that

$$n_s := \sum_{i=1}^{M} n_{i,s} \quad \text{for all strategies} s \in \{AN, AP, SN, L\}.$$

In each of the games for which there are at least $N-1$ loners, the payoff for each individual is $\sigma$. In any other case, the payoff $P_{i,s}$ of an individual of class $i$ and adopting strategy $s$ is given by

$$P_{i,L} = \sigma$$

$$P_{i,AN} = rc\frac{n_{AP}+n_{AN}}{N-n_L} - c - \alpha G^*(n_{SN})\sum_{j=1}^{M}\beta_j \cdot n_{j,AP},$$

$$P_{i,SN} = rc\frac{n_{AP}+n_{AN}}{N-n_L} - \sum_{j=1}^{M}\beta_j \cdot n_{j,AP},$$

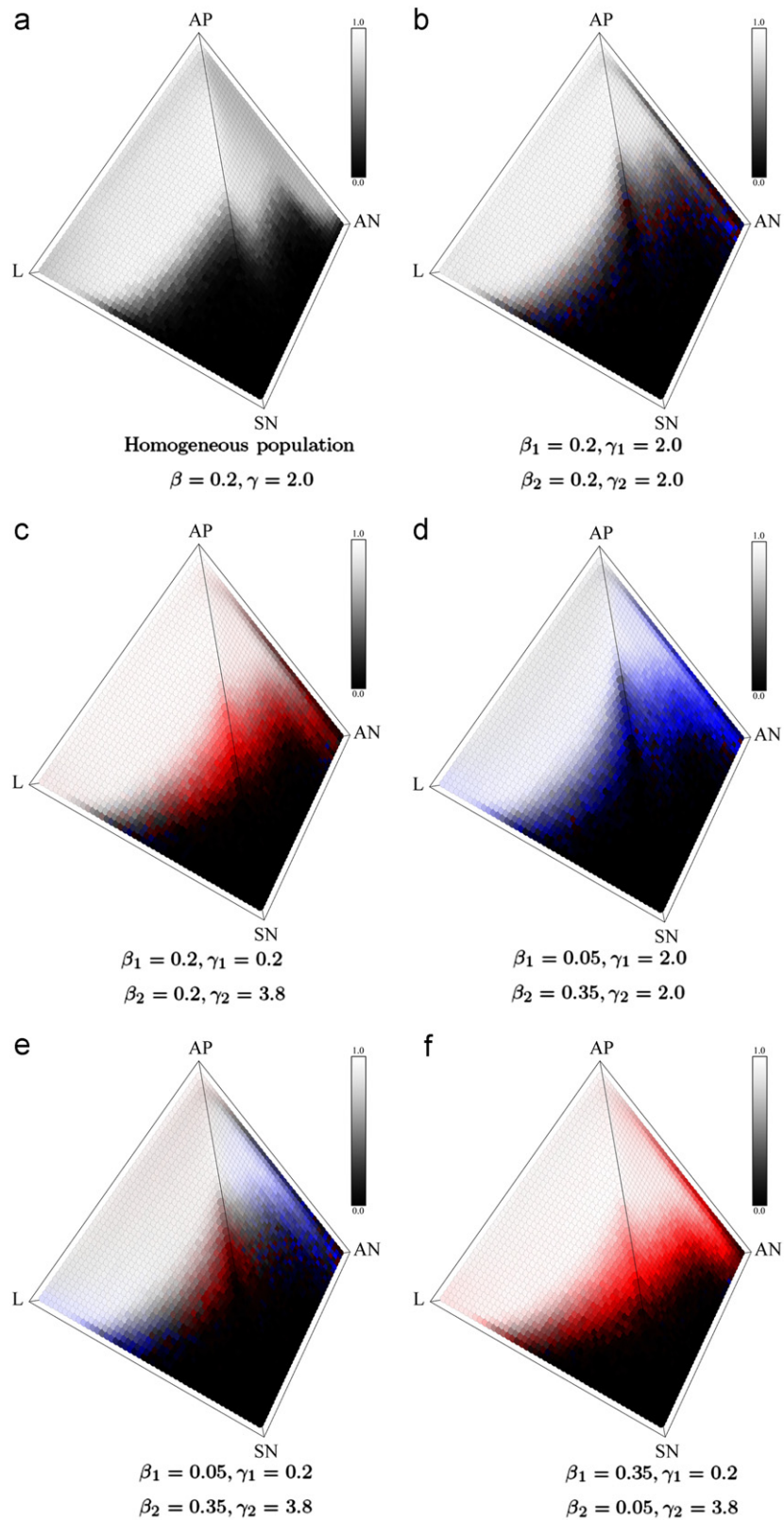$$P_{i,AP} = rc\frac{n_{AP}+n_{AN}}{N-n_L} - c - \gamma_i(n_{SN}+\alpha n_{AN}G^*(n_{SN})),$$

where $G^*(n)=1$ if $n>0$ and 0 otherwise.

### 3.2. Spatial model: results

The model described in the previous subsection has been implemented in Java in order to determine the effects of heterogeneity on the evolution of altruistic punishment in a spatially structured world. Due to the effects of assortment, the parameters for the infinite and well-mixed population of Section 2 cannot be used for the spatial model. Appendix A shows the derivation of the parameter setting for the spatial model. All results in this section are obtained by simulation on a 50 by 50 lattice[4] with periodic boundaries. Each individual simulation ran for 500 rounds of lead time, which was found to be generally sufficient for the population to reach a stable situation. Furthermore, the results were averaged over 200 simulation runs, each with a randomized initial configuration. Experiments with a larger number of repetitions did not improve the results any further.

---

[3] This type of neighborhood is also known as a Von Neumann neighborhood.
[4] Experiments with different lattice sizes showed results similar to the ones reported here.

**Fig. B.2.** Interior of the simplex $S_4$ for six different settings: (a) single homogeneous population, (b) homogeneous classes, (c) low cost, (d) low returns, (e) low returns and low cost, and (f) high returns and low cost. To improve comparability across situations, the average values $\bar{\beta}$ and $\bar{\gamma}$ over the two classes in situations (b)–(f) are fixed. The brightness of the dots indicates the proportion of altruistic individuals (AN and AP) after 500 rounds in the entire population. Color and hatching of the dots shows the relative proportion of altruistic punishers across both classes, where a red, vertically hatched dot indicates that most altruistic punishers are of class 1, while blue and horizontally hatched dots indicate that most altruistic punishers are of class 2. A more saturated color indicates a larger difference in the proportion of altruistic punishers between class 1 and 2. Proportions are averaged over 200 separate runs. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

In the spatial model, we characterize the end state of the population by the proportion of individuals that have adopted an altruistic (either punishing or non-punishing) strategy, as well as the distribution of altruistic punishers across the two classes. Fig. B.2 summarizes the results for the simulations in six different settings.[5] In every setting, the brightness of a dot represents the average proportion of altruistic individuals such that lighter dots indicate that a larger proportion of the population has adopted an altruistic strategy. The differences between classes are illustrated by the color and hatching of a dot. Red, vertically hatched dots indicate that the majority of the altruistic punishers are of class 1, while blue, horizontally hatched dots indicate that class 2 holds most of the altruistic punishers. When the number of altruistic punishers is divided equally between classes, the dot appears solid grey in Fig. B.2.

In every setting shown in Fig. B.2 the parameter values $\sigma = 0.75, c = 1.0$ and $r = 1.9$ were fixed. Furthermore, the average value of the return on punishment and the cost of punishing were constant across the simulation runs such that $\overline{\beta} = 0.2$ and $\overline{\gamma} = 2.0$. As for the results of the infinite population model, the position of a dot determines the relative frequencies of the strategies such that each corner of the simplex represents an initial state where every individual adopts the strategy indicated by the corresponding label. Note that Fig. B.2 shows the interior of the simplex, such that the proportion of each strategy is non-zero for each of the dots.

Fig. B.2a shows the results for a single homogeneous population. In contrast with the infinite population model, separating the homogeneous population into two homogeneous classes does have an effect on the results for the spatial model. In this setting, the population is homogeneous, but subdivided into two classes. The two classes interact normally, but individuals only learn behavior from individuals from the same class. Fig. B.2b shows that when individuals only learn from part of the individuals in their interaction neighborhood, altruistic behavior fails to stabilize more often in the spatially structured population. The fact that there is no inter-class learning favors selfish behavior by allowing it to exploit altruistic individuals without causing those individuals to become selfish as well. This diminishes the opportunity for positive assortment by breaking up the cluster structure that benefits cooperation and punishment in the spatial model (Nowak, 2006; Grilo and Correia, 2011).

The introduction of two homogeneous classes does not cause great variation in the number of altruistic punishers between classes. This is as expected, since there are no differences in the individual abilities between the two classes. Fig. B.2b shows mostly solid grey dots, with some more color in the dark grey area in which altruistic behavior mostly fails. In this area, altruistic punishment regularly disappears from one or both classes because of an unfavorable initial situation. The color and hatching in this area therefore indicate that one of the classes had a higher proportion of altruistic punishers over 200 separate runs by random chance.

The results when heterogeneity in the cost of punishing $\gamma_i$ is introduced, while the effectiveness of punishment $\beta_i$ remains constant across classes are shown in Fig. B.2c. The proportion of altruistic individuals appears to be fairly insensitive to changes in the cost of punishing, but heterogeneity in the cost of punishing does provide an opportunity for altruistic behavior to evolve. However, this effect is fairly limited compared to the differences in costs between classes ($\gamma_1 = 0.2$ against $\gamma_2 = 3.8$). The advantage for altruistic behavior in the case of heterogeneous cost of punishing may result in specialization between classes. As indicated by the predominantly red vertical hatching in Fig. B.2c, whenever the population does not end up in a state of all altruists or no altruists, punishment is mainly performed by class 1, the class paying the lowest cost for punishing. Note however that the effect is stronger when the initial configuration is close to the AP-SN edge, where the red and vertically hatched dots are more abundant. When the initial configuration includes more altruistic non-punishers or loners, the dots remain closer to solid grey, indicating that the distribution of altruistic punishers over the two classes is more equalized.

When instead of the cost of punishing $\gamma_i$, the returns on punishment $\beta_i$ are taken to be heterogeneous across classes, Fig. B.2d shows that heterogeneity has little effect on the evolution of cooperation. The population is slightly more likely to end up in a state with mostly altruistic individuals when the initial proportion of altruists is high, and slightly less likely to end up in such a state when selfish individuals and loners are more common in the initial configuration of strategies. However, the changes in the structure of the population are more pronounced. The majority of punishment is performed by the class of individuals with the highest returns on punishment. As opposed to the case where the cost of punishing is taken to be heterogeneous, classes that differ in the returns on punishment show the most specialization when altruists dominate the initial configuration, while the dots in Fig. B.2d remains closer to solid grey when there are many loners and selfish non-punishers.

Since heterogeneity in the cost for punishing has a positive effect on the proportion of altruistic individuals in the end state of the spatially structured population, and heterogeneity in the returns on punishment does not have a clear positive or negative effect, this leaves us with the issue of how these effects interact. Fig. B.2e shows the results when individuals of class 1 can only inflict low punishment at a low cost ($\gamma_1 = 0.2$ and $\beta_1 = 0.05$) while individuals of class 2 inflict high punishment at a high personal cost ($\gamma_2 = 3.8$ and $\beta_2 = 0.35$). Compared to the situation in which only the cost for punishing is heterogeneous, altruistic behavior has a slightly harder time to stabilize in the population. However, unlike the situation for the infinite and well-mixed population, this type of heterogeneity still represents a beneficial effect for altruistic behavior compared to the homogeneous classes.

The case of low punishment at a low cost and high punishment at a high cost also exhibits an interesting pattern of specialization between classes in the spatially structured population, combining the separate specialization effects of heterogeneity in cost for punishing and those of heterogeneity in returns on punishment. In the combined setup, which class contains most of the altruistic punishers in the final state depends on the initial configuration of strategies in the population. When altruistic non-punishers are initially rare, punishment is carried out by the class with the lowest cost, as illustrated by the red vertical hatching along the AP-SN edge in Fig. B.2e. On the other hand, when altruistic non-punishers are more common in the initial layout, the class with the highest returns on punishment ends up carrying out most of the punishment.

Finally, Fig. B.2f shows that when one group is strictly better at punishing ($\beta_1 > \beta_2$ and $\gamma_1 < \gamma_2$), this results in an advantage for altruistic behavior that is stronger than in the situation in which only the cost for punishing is heterogeneous. Moreover, the population exhibits a high degree of specialization, where almost all punishment is carried out by the efficient punishers. This specialization is mostly independent of the initial configuration of strategies. Only when there is a moderate proportion of loners in the initial configuration, the final distribution of altruistic punishers is more equalized.

---

[5] Results for the entire interior of the simplex are available online at http://www.ai.rug.nl/SocialCognition/?p=83

## 3.3. Spatial model: discussion

Where the infinite population model represents a large population in which individuals are unlikely to meet the same co-player twice in the setting of a public goods game, the spatial model described in Section 3.1 represents a relatively small population with a rigid interaction structure. Individuals only interact with a small selection of other individuals, and always in the same groups. The spatial model can therefore be used to model individuals that are arranged in a strict geographical structure, but also social connections between individuals. In the latter interpretation, the spatially structured population can be used to model the evolution and enforcing of social norms, in which the payoff of the public goods game is interpreted as a personal utility rather than evolutionary fitness.

Note that this utility is not necessarily an external reward: experiments reveal that children show a strong tendency to help others at a very young age, even in the absence of material or social rewards (Tomasello, 2009). In fact, Warneken and Tomasello (2008) show that providing a reward for helping diminishes the children's motivation to help in the future.

Social norms are customary rules of behavior that people will conform to given the expectation that others will conform to it too (Lewis, 1969; Bicchieri, 2006). The individuals in our model have no way of explicitly forming expectations about the actions of others. They simply imitate the strategy of the individual with the highest payoff. The implicit expectation is that imitating the most successful strategy will raise their own payoff. In the setting of the public goods game, altruistic strategies yield the highest payoff when they are used by everyone else in the interaction neighborhood. Altruistic punishment can therefore become a social norm, since individuals prefer to imitate the altruistic behavior on the condition that everyone else is an altruistic punisher.

As expected from results of two-player games (Nowak and May, 1993; Lindgren and Nordahl, 1994; Killingback and Doebeli, 1996), altruistic behavior and punishment in the public goods game can evolve more easily in a spatial model than in the infinite and well-mixed population. Because of the localized interactions, altruists are at a lower risk of being exploited by selfish individuals, while punishers only punish selfish individuals which they will encounter again in the public goods game. Moreover, if selfish behavior is highly profitable, some of the altruistic individuals that the selfish individual has interacted with will imitate the strategy, withholding cooperation the next round. In this sense, the spatial structure provides a form of direct reciprocity.

We have shown that in a spatially structured population, heterogeneity in the returns on punishment by itself has little effect on whether or not altruistic behavior will stabilize. Differentiation in the cost of punishing does have a positive effect for altruistic behavior, although the effect does not seem to be as pronounced as in the case of the infinite and well-mixed population. However, heterogeneous classes readily specialize in the spatial model such that punishment is carried out by the class of individuals best suited for the task, whether this is because they enjoy a higher return on punishment or because they pay a lower cost for punishing. When a class of individuals combines both benefits to their ability to punish, inflicting a higher punishment at a lower cost, punishment becomes the almost sole responsibility of this class. Compared to the homogeneous case, heterogeneity of this kind has a clear positive effect on the evolution of altruistic behavior in the sense that fewer punishers are needed in the initial configuration to reliably stabilize altruistic behavior in the entire population. As with the infinite and well-mixed population, the combined effect of heterogeneity in the costs of punishing and the returns on punishment is stronger than the separate effects.

The ability of the spatial model to specialize is best illustrated in the case where the two heterogeneous classes differ in their level of punishment, such that one class inflicts high punishment at a high cost, while the other inflicts a lower punishment at a lower cost. In this case, which class will specialize into becoming the class of punishing individuals depends on the initial proportions of strategies. When altruistic behavior is already common, the class that inflicts high punishment will contain most altruistic punishers. On the other hand, if altruistic behavior is less common, and punishers are faced with more selfish individuals to punish, the class enjoying a lower cost for punishing will take over the responsibility of punishing.

In the interpretation of social norms that are maintained by a group of people, the spatial model predicts that punishment will be carried out either by the ones paying the lowest personal cost or the ones inflicting the highest punishment. When there is a choice between low punishment at a low cost and high punishment at a high cost, the level of punishment depends on the popularity of the social norm. When a social norm is popular in the sense that many individuals adhere to it, the model predicts that violation of the norm will be punished severely, even if it comes at a high personal cost to the punishers. When only few individuals adhere to the norm in the initial situation, punishment for violating it will be lower. Note however that this choice in the level of punishment is not made individually in our model. Individuals may choose whether or not they punish for selfish behavior, while the level of punishment is entirely determined by their individual abilities.

The specialization that occurs in the spatially structured population shows some similarities with experimental results on the bystander effect. In general, bystanders are slower to help and help less often during an emergency situation when other bystanders are present (Latane and Darley, 1968). However, experimental research has shown that this bystander effect does not occur when subjects consider other bystanders to be unable to help (Bickman, 1971). Moreover, when subjects consider themselves to be more competent in dealing with the emergency, the presence of other bystanders does not inhibit helping either (Pantin and Carver, 1982; Cramer et al., 1988). In naturally occurring situations, bystanders that intervene in a crime generally describe themselves as being physically strong and appear to act out of a sense of capability through training experience or personal strength (Huston et al., 1981). On the other hand, preschoolers are less likely to respond to the distress of one of their peers when a competent adult caregiver is present (Caplan and Hay, 1989). This sense of responsibility, where emergencies and norm violations are handled by the individual most competent to complete the task, also appears in the spatially structured population. Our results show that even if the individual abilities of others are not observable, norm violation will be punished by the individuals best suited for the task. However, if none of these individuals are present in the group playing the public goods game, norm violation is likely to go unpunished.

## 4. Conclusions and future research

Most models in the literature concerned with the question of how altruistic and punishing behavior may have evolved assume that the population consists of homogeneous individuals (Nowak and May, 1993; Lindgren and Nordahl, 1994; Killingback and Doebeli, 1996; Hauert et al., 2009; Eldakar et al., 2007; Eldakar and Wilson, 2008). Each individual in the population has the same costs and benefits of cooperating, as well as the same costs and effectiveness of punishing for selfish behavior. We have shown that relaxing this assumption of homogeneity by allowing for

populations that consist of two or more different sub-classes affects the evolution of altruistic punishment. In this section, we will present our conclusions and provide directions for future research.

### 4.1. Conclusions

In this paper, we investigated the effects of heterogeneity on the evolution of altruistic punishment. In particular, we determined how individual differences in the cost inflicted by punishment and the personal cost at which punishment may be performed influence the conditions under which altruistic behavior can stabilize in a population, as well as the resulting structure of altruistic punishers in the population. To achieve this, we extended the public goods game with voluntary participation of Brandt et al. (2006) by separating the population into two classes such that individuals within the same class are homogeneous, but may differ in the cost they pay to punish for selfishness and the cost punishers inflict on selfish individuals between classes. Furthermore, we compared these results for two population models: an infinite size population and a spatially structured population.

Based on these results, heterogeneity of individuals may certainly have an effect on the evolution of altruistic punishment. In the setting of an *infinite and well-mixed* population, individuals that meet in a public goods game are unlikely to meet each other again in their lifetime. Heterogeneity in the cost for punishing in this setting can make it easier for altruistic behavior to stabilize; compared to a homogeneous population in which the average cost for punishing is the same, a lower proportion of punishers is needed to ensure that the population ends up in a state with only altruistic individuals. In contrast, introducing individual differences in the effectiveness of punishment has no influence on whether or not altruistic behavior will stabilize in an infinite population by itself. However, when combined with heterogeneity in the cost of punishing such that individuals that inflict high punishment do so at a lower cost, individual differences in the effectiveness of punishment can amplify the positive effect on the evolution of altruistic behavior. For the infinite and well-mixed population, the interaction between heterogeneity of cost for punishing and heterogeneity of returns on punishment can work both ways. When individuals differ in the level of punishment, such that some individuals inflict a low punishment at a low personal cost, while others inflict high punishment at a high personal cost, altruistic behavior can become harder to stabilize.

In addition to the infinite and well-mixed population, we also investigated the effects of heterogeneity in the individual abilities to punish in a *spatially structured* population. In this setting, individuals are assigned a spatial location on a grid, and only play the public goods game within a local neighborhood. As in the infinite population model, we found a positive effect of heterogeneity in the cost for punishing. This effect was amplified when individuals that inflict high punishment pay a low cost to do so. In contrast to the infinite population model, we found heterogeneity in the returns of punishment makes it slightly easier to stabilize altruistic behavior when it is common in the initial situation, but slightly harder when altruistic behavior is initially rare. This result becomes even more clear when individuals differ in their level of punishment. When altruistic behavior is initially common, individuals that inflict high punishment at a high cost will perform most of the punishment. However, when altruistic behavior is initially rare, punishment becomes the responsibility of individuals that inflict low punishment at a low cost.

Note that in our analysis, individual differences in the returns on and costs for punishing do not affect the share of the public good the punisher is entitled to. An individual that can punish with unusually high effectiveness gains no more benefits from punishing than others: the increased effect of discouraging selfish behavior is of no greater advantage to them than to others. In nature, dominant animals in a social group are more likely to punish for failure to contribute to the public good, but they commonly gain a disproportionate share of reproduction. This is not simply the right of the strongest, since subordinates are known to challenge a dominant animal if it takes more than its fair share (Clutton-Brock and Parker, 1995). This additional immediate benefit of punishing for selfish behavior ensures that the punishing individuals have higher stakes in stabilizing cooperation in the population.

In our model such additional advantages do not exist. The effects of heterogeneous effectiveness of punishment and the resulting structure of the population are purely caused by differences in the abilities to punish. Based on these differences, individuals may pay a lower personal cost for punishing, but the act of punishing always represents a short-term loss in fitness for the punisher.

As a final note, in the models presented here we assumed that there is no mutation. Due to the rock-paper-scissors dynamics of loners, altruistic non-punishers and selfish non-punishers, adding mutation in the spatial model prevents convergence by re-introducing strategies that have disappeared from the population (see also Hauert et al., 2007). Appendix B shows results of additional experiments that introduce mutation in the spatial model. Although the initial configuration of strategies no longer affects the eventual fate of the population, the results are similar to the ones described in Section 3.

### 4.2. Future research

A number of issues are left open for further research. In our models, we assumed that the population consisted of two separate classes of equal size, where the heterogeneity was limited between classes, while individuals within the same class were homogeneous. We also assumed that individuals were aware of this division in classes, and were able to determine the class of individuals they interacted with. In general, however, individuals may have continuously distributed abilities, which others may not be able to determine correctly.

In the models we discussed, we allowed for only one type of punishment. Altruistic punishers impose a fee on selfish individuals, as well as on altruists that fail to punish. However, selfish individuals also have incentive to discourage selfishness of others, which leads to selfish punishment (Nakamaru and Iwasa, 2006; Eldakar et al., 2007; Eldakar and Wilson, 2008). This type of punishment also appears in human (Falk et al., 2005) as well as in animal societies (Wenseleers et al., 2005). Furthermore, punishment can also be directed at loners (Hauert et al., 2007) or altruistic individuals (Nikiforakis and Normann, 2008; Herrmann et al., 2008; Gächter and Herrmann, 2009), which may have a strong negative effect on the evolution of cooperative behavior (Herrmann et al., 2008; Janssen and Bushman, 2008; Dreber et al., 2008; Gächter et al., 2008; Wu et al., 2009; Rand et al., 2010). Finally, rather than a system of peer-punishment, in which punishers individually punish for selfishness and are revealed only after playing the public goods game, a system of pool-punishment is also a possibility (Sigmund et al., 2010). It remains an open question how different types of punishing interact with the effects of heterogeneity in the effectiveness and costs of punishment.

The differences between the infinite and well-mixed population and the spatially structured population are quite large. In the infinite population model there are no repeat encounters. A pair of individuals that meet in the setting of the public goods game never meet each other in the same setting again. On the other hand, individuals in the spatially structured population only interact with the same 12 neighbors, playing the same five public

goods games every round. However, the public goods game is well suited for an intermediate form, in which co-players are selected at random from a local neighborhood, combining elements of both population models.

The simplified representation of individuals raises another issue. Each individual adopts one strategy and uses that strategy in all the games it plays. Individuals cannot take advantage of past experience when they encounter a co-player they have met before, which is particularly relevant in the spatially structured population. This simplified representation precludes the emergence of the hierarchical structure that is commonly found in animal societies. Research on dominance relations in competitive environments has resulted in models that allow for these hierarchical structures, such as DomWorld (Hemelrijk, 1999, 2000, 2002). Future research could shed light on how these competitive hierarchical structures affect cooperative efforts.

In our model of the public goods game, altruistic behavior is represented as a binary choice; individuals either invest in the public good, or not. In practice, the amount invested in the public good may be chosen from a continuous range of possibilities, depending on individual abilities. Empirical research has shown that human subjects readily accept individual differences and adjust their expectations, punishing only when they believe their co-players invested less than their fair share (Fisher et al., 1995; Reuben and Riedl, 2009). In our models, we have shown that a sense of "fairness" does seem to emerge in the spatial model, in the sense that punishment is carried out by individuals that are best suited for the task. This specialization occurs even in the absence of personal benefits. Whether this primitive form of a "fairness" principle also emerges in a more complex setting where cooperative behavior is not simply a binary choice is a question for future research.

In conclusion, it turns out that heterogeneity provides additional opportunities for cooperative behavior to be maintained and spread through a population of individuals. But many intriguing questions remain on how this heterogeneity exactly evolves and what role it plays in animal societies.

## Acknowledgments

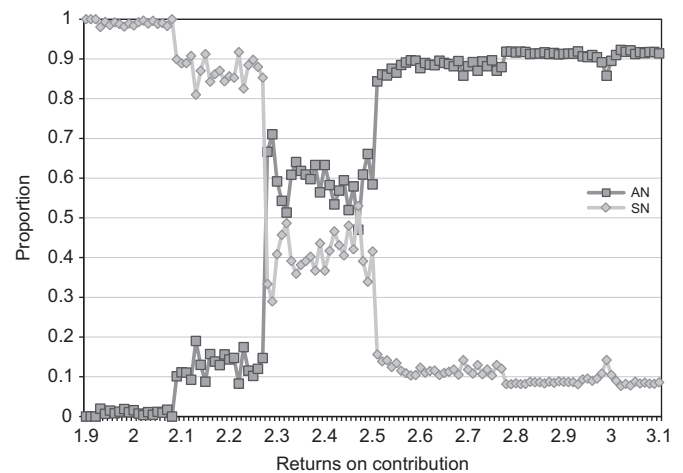## Appendix A. Parameters of the spatial model

It is known from game-theoretical games such as the prisoner's dilemma that introducing spatial structure into the model of a population can have a powerful effect on the behavior of the individuals. Even without the help of punishing individuals, small clusters of altruistic non-punishers already provide enough of an advantage to prevent them from copying selfish behavior (see among others Nowak and May, 1993; Lindgren and Nordahl, 1994; Killingback and Doebeli, 1996). The basic public goods model we use as a starting point is no exception. Because of this, the parameters used for the infinite and well-mixed population of Section 2 do not produce comparable results when implemented for the spatial model. In this section, we therefore describe the parameter setting for the spatial model.

One of the effects of spatial structure is that due to localized interactions, selfish individuals can no longer exploit distant altruistic individuals, which encourages altruistic behavior. Fig. B.3 shows this result by showing the proportion of altruistic and selfish
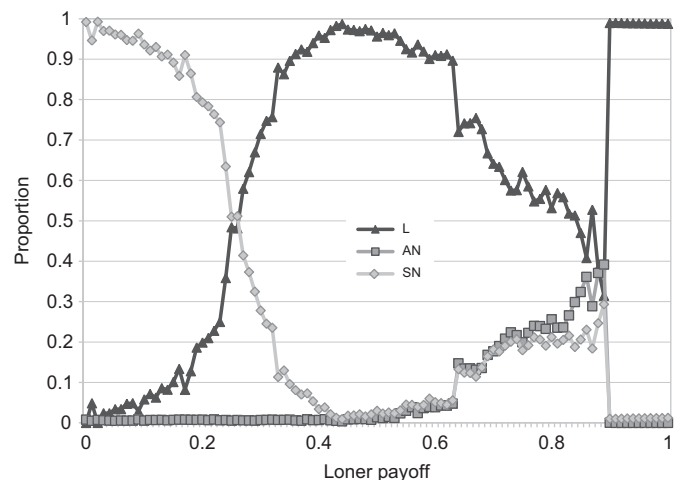
individuals for different levels of the return on contribution $r$. In this setting, every individual is randomly assigned the strategy of altruistic or selfish non-punisher, after which the proportions of the two strategies after a 500 round lead time was determined. Fig. B.3 shows these proportions averaged over 200 runs, for every 0.01 change of $r$ in the range $1.9 \leq r \leq 3.1$.

In the infinite population model, altruistic non-punishers cannot survive in the absence of loners and punishers, no matter how high the return on contribution. Fig. B.3 shows that in the lattice model, sufficiently dense clusters of altruistic non-punishers can withstand an invasion of selfish non-punishers for $r > 2.1$. For $r > 2.28$, altruistic non-punishers can even outperform selfish non-punishers in terms of fitness, and will reliably represent over 80% of the population for $r > 2.5$. For the purposes of comparing the effects of heterogeneity on the behavior of a population on a lattice, we therefore set the return on contribution $r=1.9$ instead of using $r=3.0$ as in the setting of an infinite size population.

To make sure that the public good remains competitive, the loner payoff $\sigma$ should be lower than the maximum payoff for altruistic individuals $(r-1)c$. For $r=1.9$ and $c=1.0$, Fig. B.4 shows



**Fig. B.3.** Effect of the value of returns on contribution ($r$) on the proportions of altruistic non-punishers (AN) and selfish non-punishers (SN). The initial distribution of strategies was randomized for each run. The final proportions were determined after 500 rounds of lead time and averaged over 200 runs, for every 0.01 change of $r$ in the range $1.9 \leq r \leq 3.1$.



**Fig. B.4.** Effect of the value of loner payoff ($\sigma$) on the proportions of loners (L), altruistic non-punishers (AN) and selfish non-punishers (SN). The initial distribution of strategies was randomized for each run. The final proportions were determined after 500 rounds of lead time and averaged over 200 runs, for every 0.01 change of $\sigma$ in the range $0 < \sigma \leq 1.0$.
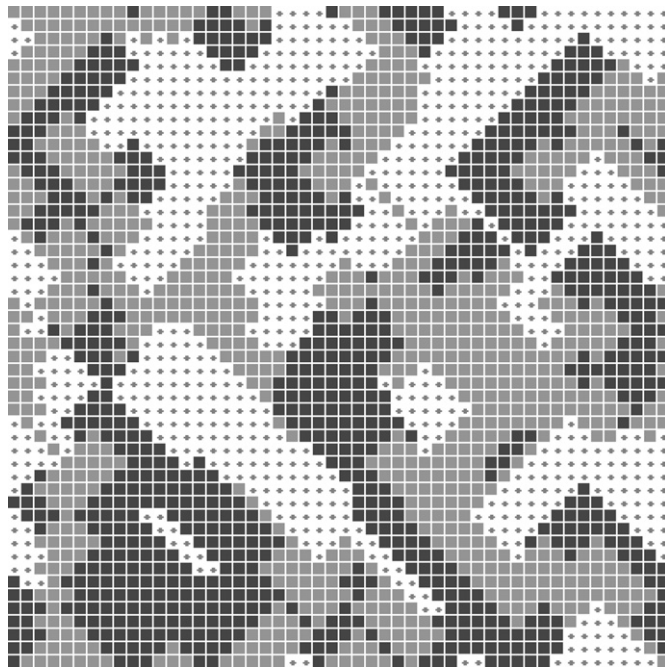
the proportions of loners, altruistic non-punishers and selfish non-punishers as a function of the loner payoff $\sigma$. As before, the results are averaged over 200 runs. In each of these runs, the population was randomly initialized and given 500 rounds of lead time before the proportions of the three strategies were determined. This process was repeated for every 0.01 change in $\sigma$ in the range $0 < \sigma \leq 1.0$. Note that the value $\sigma = 0$ was omitted for display purposes.

For $\sigma \geq 0.9$, altruistic behavior is at a disadvantage to a loner, which means that the population will eventually end up in the situation of 100% loners. Although some individuals may adopt a



**Fig. B.5.** Typical results of a population with $r = 2.0$, $c = 1.0$ and $\sigma = 0.8$. Selfish non-punishers (dark grey) exploit altruistic non-punishers (light grey), until their payoff decreases below the loner (white) payoff.
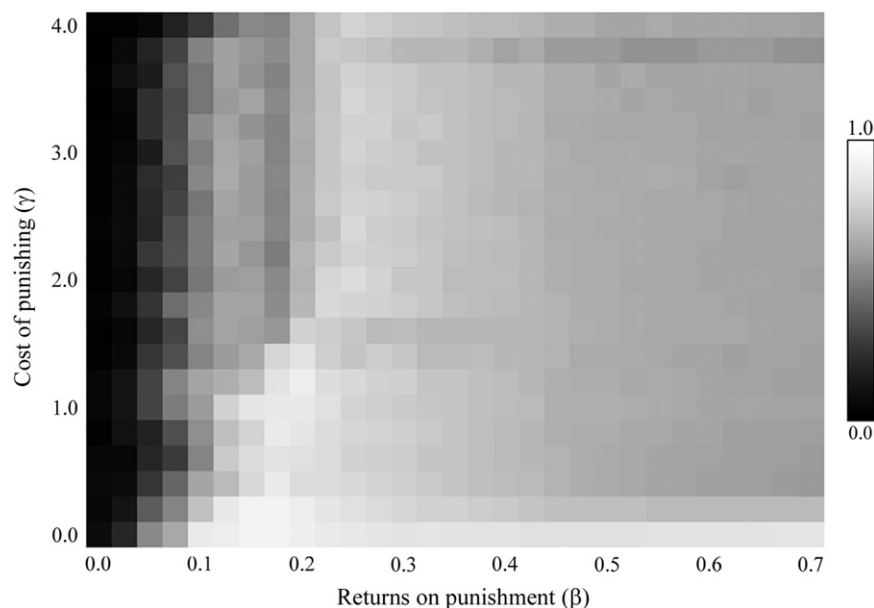
non-loner strategy in this situation, none of them will have any other non-loner individual in their interaction neighborhood, and they are therefore forced to play as a loner as well.

Another discrete step in the graph appears around $\sigma = 0.63$. To avoid artefacts caused by either one of these steps, a loner payoff of $\sigma = 0.75$ was chosen for the public goods game simulations. Fig. B.5 shows a typical result of this parameter setting when there are no punishers. As in the case with the infinite population model, altruistic individuals, selfish individuals and loners are locked in an infinite game of rock-paper-scissors (Nowak, 2006). Selfish individuals exploit the altruistic non-punishers, but the popularity of their strategy quickly causes individuals in their interaction neighborhood to adopt the selfish non-punisher strategy as well, sharply reducing their payoff. This leaves all selfish non-punishers open to become a loner. In the absence of selfishness, the payoff of altruistic non-punishers increases, tempting the loners to rejoin the public goods game. Unlike the infinite population model, in which the initial configuration of strategies determines which periodic orbit the population will eventually reach, the initial state of the population has little effect on the eventual proportions of strategies in the lattice model. Typically, in the first rounds after initialization of the population, most of the altruistic and selfish individuals are replaced by loners, after which remaining clusters of altruists start expanding and the situation of Fig. B.5 appears.

It is well known from research on two-person games such as the prisoner's dilemma (Nowak and May, 1993) and hawks and doves (Killingback and Doebeli, 1996) that punishing is much more efficient on a lattice than in an infinite and well-mixed population. Because of the local interactions, once a cluster of altruistic punishers has appeared, it can easily grow. The interior of the cluster contains no selfish individuals, which means that the altruistic punishers enjoy the full benefit of their mutual cooperation, without the burden of having to punish for selfishness. Meanwhile, on the edge of the cluster, selfish behavior is punished, causing the payoff for selfishness to decrease. Even though punishment is costly, as long as the high payoff of other altruistic punishers in the interior of the cluster is higher than the payoff of the punished selfish individuals, the individuals on the edge will not change their strategy.
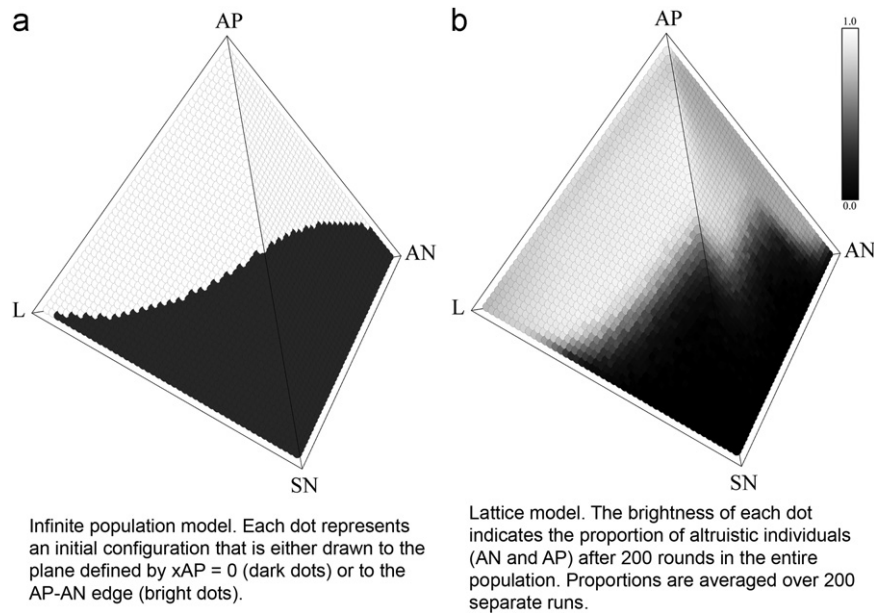


**Fig. B.6.** Effect of the values of the returns on punishment $\beta$ and cost of punishing $\gamma$ on the proportion of altruistic punishers. The initial distribution of strategies was randomized for each run, such that AP and SN represent 40% of the population each, and AN and L represent 10% each. The final proportions were determined after 500 rounds of lead time and averaged over 200 separate runs, for every 0.025 change of $\beta$ in the range $0 \leq \beta \leq 0.7$ and every 0.2 change of $\gamma$ in the range $0 \leq \gamma \leq 4.0$.

Infinite population model. Each dot represents
an initial configuration that is either drawn to the
plane defined by xAP = 0 (dark dots) or to the
AP-AN edge (bright dots).

Lattice model. The brightness of each dot
indicates the proportion of altruistic individuals
(AN and AP) after 200 rounds in the entire
population. Proportions are averaged over 200
separate runs.

**Fig. B.7.** Interior of the simplex $S_4$ for (a) the infinite and well-mixed population model and (b) the lattice model, both for a homogeneous population.

Fig. B.6 shows that this effect also holds for the public goods game. The figure shows the proportion of altruistic punishers after 500 rounds of the public goods game as a function of the return on punishment $\beta$ and the cost for punishing $\gamma$. In this setting, populations were initialized on a 50 by 50 lattice such that approximately 40% of the population started as altruistic punisher, 40% started as selfish non-punisher, and the remaining 20% consisted of loners and altruistic non-punishers. The proportion of altruistic punishers was recorded after 500 rounds of play, and was averaged over 200 separate runs. Only when the return on punishment fell short of 0.2 did the proportion of altruistic punishers drop below 50%.

As shown by Fig. B.6, the evolution of altruistic punishment is largely insensitive of the cost of punishing. The reason behind this is that due to the synchronous updating, altruistic individuals expand in clusters. Altruistic punishers at the edge of such a cluster are forced to punish many selfish individuals outside of the cluster, sharply reducing their fitness. However, punishers need not reduce the fitness of selfish individuals below that of their own. When the fitness of the selfish individuals is lower than the fitness of any altruistic punisher in the cluster, where there are few selfish individuals to punish, altruistic punishers at the edge of the cluster will not change their strategy.

The parameters used for the lattice model and a homogeneous population are $N=5, \sigma=0.75, c=1.0, r=1.9, \alpha=0.1, \gamma=2.0$ and $\beta=0.2$. Note that these parameters are less favorable for altruistic individuals and punishers when compared to the parameter setting for the infinite and well-mixed population of Section 2.2 ($N=5, c=\sigma=\gamma=1.0, r=3.0, \beta=1.2, \alpha=0.1$). Fig. B.7 shows the results of a homogeneous population for the infinite and well-mixed population model and the lattice model side by side. Fig. B.7a is repeated from Section 2.2. In this figure, bright dots indicate initial situations that eventually end up in a state in which every individual is altruistic, whereas dark dots indicate that such a situation will never occur. Fig. B.7b shows the results for a lattice model, where the brightness of each dot indicates the proportion of individuals that are altruistic after 500 rounds of playing the public goods game, where brighter dots indicate more altruistic individuals in the population. For the lattice model, the results are averaged over 200 separate runs. Even though the parameters are less favorable for the lattice model, the results
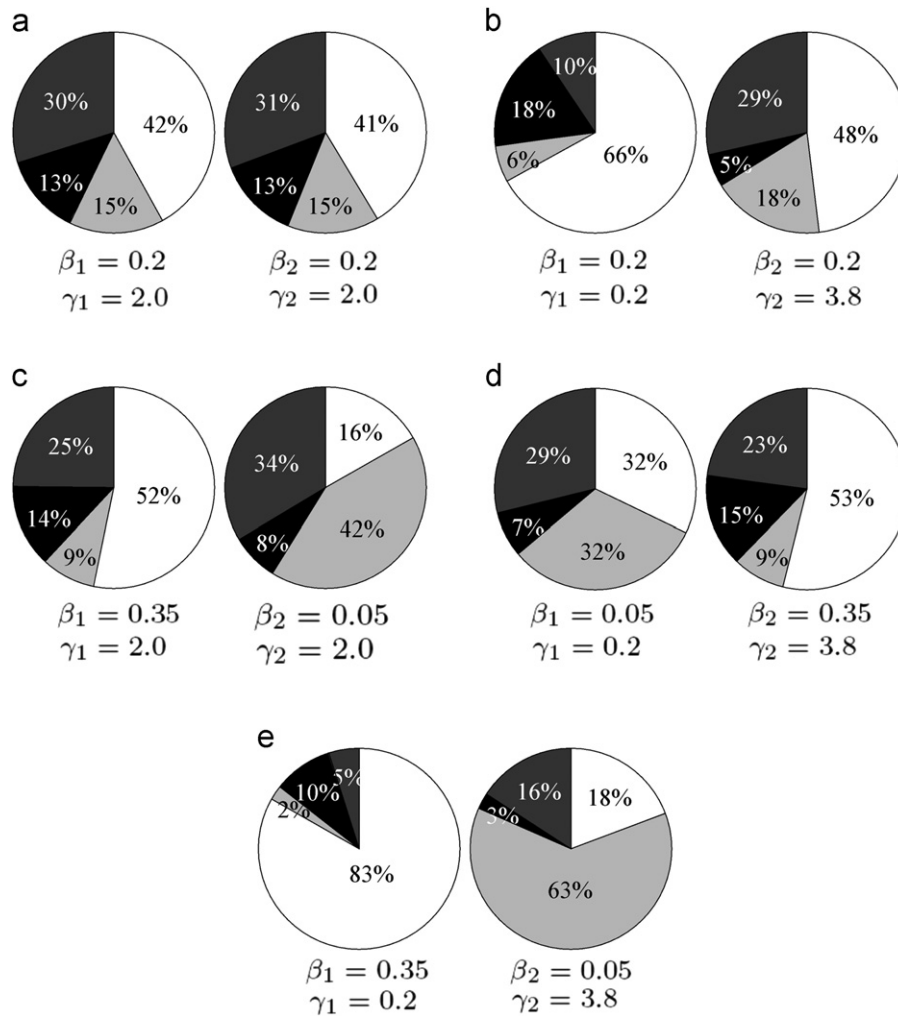
compare reasonably well to the results for the infinite and well-mixed population model. This parameter setting is therefore used as the base scenario for determining the effects of heterogeneity in Section 3.2.

## Appendix B. Effects of mutation on the spatial model

In the models presented in Section 3, we assumed that no mistakes are made in determining which individual has the highest payoff and in performing the actions associated with the strategy an individual has adopted. However, noise and mistakes can greatly influence the evolution of altruistic behavior (see among others Foster and Young, 1990; Fudenberg and Maskin, 1990; Fudenberg and Harris, 1992; Blume, 2003), especially in spatial games. In this section, we therefore examine the effects of mutation on our results for the spatially structured population of Section 3.

We implemented mutation as a small probability that an individual adopts a randomly chosen strategy rather than choosing the one in its learning neighborhood that received the highest payoff. This way, mutation allows the reintroduction of strategies that had previously been eliminated from the population. Due to the rock-paper-scissors dynamics of selfish non-punishers, loners and altruistic non-punishers, mutation therefore effectively prevents the population from converging into a situation where all individuals share the same strategy of either non-punishers, loners or altruistic non-punishers. Furthermore, in the absence of selfish individuals, second-order free-riding cannot be detected, which means that through neutral drift, a population consisting only of altruistic punishers can be invaded by altruistic non-punishers. In Hauert et al. (2007), the population is therefore modeled using a Moran process to derive the average time a system spends in each absorbing state where all individuals share the same strategy. In this section, we will determine the effects of heterogeneity in the effectiveness and costs of punishing on these durations using agent-based simulation.

Fig. B.8 shows the results of our experiments. In each situation, the population is separated into two classes, such that individuals within each class are homogeneous, but individuals may differ in their ability to punish across classes. For each class, the pie charts show the average amount of time an individual spends using each

**Fig. B.8.** Proportion of altruistic punishers (white), altruistic non-punishers (light gray), loners (dark gray) and selfish non-punishers (black) in each class after 20,000 rounds. Results were averaged over 500 runs.

of the four possible strategies selfish non-punisher (black), loner (dark gray), altruistic non-punisher (light gray) and altruistic punisher (white). These results were obtained by measuring the proportion of individuals adopting each strategy after 20,000 rounds in the game, averaged over 500 runs. For each run, the population was randomly initialized with an equal number of individuals per class, who initially adopted a randomly drawn strategy. The game was played on a 50 by 50 lattice[6] with periodic boundaries. The rate of mutation has been chosen such that on average, one individual per time unit mutates ($\mu = 0.0004$).

Fig. B.8a shows the results for a population that is separated into two homogeneous groups. As expected, individuals from each class spend the same amount of time as an altruistic punishers (41%), while they are altruistic 56% of the time.

Fig. B.8b shows that when classes differ in the cost they pay for punishing, individuals spend more time as an altruistic individual. Individuals that pay a low cost for punishing are altruistic 72% of the time, whereas the other class behaves altruistically 66% of the time. The responsibility for punishing for selfishness also shifts toward the class that pays a lower cost for punishing. However, although these individuals spend more time punishing for self-ishness, they also spend more time behaving selfishly themselves.

When instead of the cost for punishing the effectiveness of punishment is heterogeneous, this has a less pronounced effect on altruistic behavior, as shown in Fig. B.8c. Individuals of either class spend approximately 60% of their time being altruistic. However, the heterogeneity in effectiveness of punishment changes the distribution of punishers across classes. Individuals of the class that is more effective in punishing spends over three times as much time punishing than those that are less effective. In this case, the average time spent as a punisher across classes is less than the base situation shown in Fig. B.8a.

When the two types of heterogeneity are combined, such that one class pays a low cost ($\gamma_1 = 0.2$) to impose low punishment ($\beta_1 = 0.05$) for selfishness, while individuals from the other class pays a high cost ($\gamma_2 = 3.8$) to impose a higher punishment ($\beta_2 = 0.35$), the proportion of punishers in the population increases compared to the situation that only includes heterogeneity in the effectiveness of punishment. Although they pay a high cost to do so, individuals that impose high punishment spend over half their time (53%) punishing others. However, the low cost of the lower punishment also encourages individuals from the other class to spend more time as a punisher (32%) compared to the situation where only effectiveness of punishment is varied (17%). The additional punishers, increase overall altruistic behavior such that on average, individuals spend 63% of their time behaving altruistically.

When high punishment can be achieved at a low cost for one of the classes, while individuals in the other class pay a higher

---

[6] Our experiments show that although larger lattices tend to favor altruistic behavior, the effects of heterogeneity are similar to those reported here.

cost to impose a lower fine, the effects on altruistic behavior are more dramatic. Fig. B.8e shows that across classes, individuals spend 83% of their time altruistically, with individuals that punish efficiently spending almost all that time punishing others. Efficient punishers also spends four times as much time punishing others for selfish behavior as inefficient punishers. However, individuals in the class of efficient punishers also spend three times as long being selfish as individuals from the other class.

In this section, we have shown how adding mutation to the spatial model presented in Section 3 affects the results we obtained. We find that similar to the situation without mutation, heterogeneity in the effectiveness of punishment and costs for punishing increases the proportion of time individuals spend adopting an altruistic strategy across both classes. The effects on the structure of the population are similar as well. In general, individuals that are better suited to be punishers, either because they pay a lower cost to punish, or because they incur a higher punishment, spend more of their time as an altruistic punisher. Altruistic behavior benefits most from heterogeneity that separates the population into efficient and ineffective punishers, such that one class pays a lower cost to incur higher punishment than the rest of the population. In this case, the time that individuals spend altruistically is highest, while efficient punishers are responsible for most of the punishment.

## References

Axelrod, R., Hamilton, W., 1981. The evolution of cooperation. Science 211 (4489), 13–90.

Bicchieri, C., 2006. The Grammar of Society. Cambridge University Press, Cambridge, UK.

Bickman, L., 1971. The effect of another bystander's ability to help on bystander intervention in an emergency. J. Exp. Soc. Psychol. 7 (3), 367–379.

Blume, L., 2003. How noise matters. Games Econ. Behav. 44 (2), 251–271.

Bolton, G., Zwick, R., 1995. Anonymity versus punishment in ultimatum bargaining. Games Econ. Behav. 10 (1), 95–121.

Boyd, R., Gintis, H., Bowles, S., Richerson, P., 2003. The evolution of altruistic punishment. Proc. Natl. Acad. Sci. 100 (6), 3531–3535.

Brandt, H., Hauert, C., Sigmund, K., 2003. Punishment and reputation in spatial public goods games. Proc. R. Soc. Lond. Ser. B: Biol. Sci. 270 (1519), 1099–1104.

Brandt, H., Hauert, C., Sigmund, K., 2006. Punishing and abstaining for public goods. Proc. Natl. Acad. Sci. 103 (2), 495–497.

Camerer, C., Thaler, R., 1995. Anomalies: ultimatums, dictators and manners. J. Econ. Perspect. 9 (2), 209–219.

Caplan, M., Hay, D., 1989. Preschoolers' responses to peers' distress and beliefs about bystander intervention. J. Child Psychol. Psychiatry 30 (2), 231–242.

Clutton-Brock, T., Parker, G., 1995. Punishment in animal societies. Nature 373 (6511), 209–216.

Cramer, R., Mcmaster, M., Bartell, P., Dragna, M., 1988. Subject competence and minimization of the bystander effect. J. Appl. Soc. Psychol. 18 (13), 1133–1148.

Crespi, B., 2001. The evolution of social behavior in microorganisms. Trends Ecol. Evol. 16 (4), 178–183.

Dreber, A., Rand, D., Fudenberg, D., Nowak, M., 2008. Winners don't punish. Nature 452 (7185), 348–351.

Dugatkin, L., 1997. Cooperation Among Animals: An Evolutionary Perspective. Oxford University Press, Oxford, UK.

Eldakar, O., Wilson, D., 2008. Selfishness as second-order altruism. Proc. Natl. Acad. Sci. 105 (19), 6982–6986.

Eldakar, O., Farrell, D., Wilson, D., 2007. Selfish punishment: altruism can be maintained by competition among cheaters. J. Theor. Biol. 249 (2), 198–205.

Falk, A., Fehr, E., Fischbacher, U., 2005. Driving forces behind informal sanctions. Econometrica, 2017–2030.

Fehr, E., Gächter, S., 2000. Cooperation and punishment in public goods experiments. Am. Econ. Rev., 980–994.

Fehr, E., Gächter, S., 2002. Altruistic punishment in humans. Nature 415 (6868), 137–140.

Fisher, J., Isaac, R., Schatzberg, J., Walker, J., 1995. Heterogenous demand for public goods: behavior in the voluntary contributions mechanism. Public Choice 85 (3), 249–266.

Fletcher, J., Zwick, M., 2007. The evolution of altruism: game theory in multilevel selection and inclusive fitness. J. Theor. Biol. 245 (1), 26–36.

Foster, D., Young, P., 1990. Stochastic evolutionary game dynamics. Theor. Popul. Biol. 38 (2), 219–232.

Fowler, J., 2005. Altruistic punishment and the origin of cooperation. Proc. Natl. Acad. Sci. 102 (19), 7047–7049.

Fudenberg, D., Harris, C., 1992. Evolutionary dynamics with aggregate shocks. J. Econ. Theor. 57 (2), 420–441.

Fudenberg, D., Maskin, E., 1990. Evolution and cooperation in noisy repeated games. Am. Econ. Rev. 80 (2), 274–279.

Gärdenfors, P., 2011. The cognitive and communicative demands on cooperation. In: van Eijck, J., Verbrugge, R. (Eds.), Games, Actions, and Social Software-Springer, Berlin.

Güth, W., Schmittberger, R., Schwarze, B., 1982. An experimental analysis of ultimatum bargaining. J. Econ. Behav. Organ. 3 (4), 367–388.

Griffin, A., West, S., 2002. Kin selection: fact and fiction. Trends Ecol. Evol. 17 (1), 15–21.

Grilo, C., Correia, L., 2011. Effects of asynchronism on evolutionary games. J. Theor. Biol. 269 (1), 109–122.

Gächter, S., Herrmann, B., 2009. Reciprocity, culture and human cooperation: previous insights and a new cross-cultural experiment. Philos. Trans. R. Soc. B: Biol. Sci. 364 (1518), 791.

Gächter, S., Renner, E., Sefton, M., 2008. The long-run benefits of punishment. Science 322 (5907), 1510.

Hamilton, W., 1964. The genetical evolution of social behaviour. II. J. Theor. Biol. 7 (1), 17–52.

Hardin, G., 1968. The tragedy of the commons. Science 162 (859), 1243–1248.

Hauert, C., DeMonte, S., Hofbauer, J., Sigmund, K., 2002. Volunteering as red queen mechanism for cooperation in public goods games. Science 296 (5570), 1129–1132.

Hauert, C., DeMonte, S., Hofbauer, J., Sigmund, K., 2002. Replicator dynamics for optional public good games. J. Theor. Biol. 218 (2), 187–194.

Hauert, C., Traulsen, A., Brandt, H., Nowak, M., Sigmund, K., 2007. Via freedom to coercion: the emergence of costly punishment. Science 316 (5833), 1905–1907.

Hauert, C., Traulsen, A., Brandt, H., Nowak, M., Sigmund, K., 2009. Public goods with punishment and abstaining in finite and infinite populations. Biol. Theor. 3 (2), 114–122.

Hemelrijk, C., 1999. An individual-orientated model of the emergence of despotic and egalitarian societies. Proc. R. Soc. Lond. Ser. B: Biol. Sci. 266 (1417), 361.

Hemelrijk, C., 2000. Towards the integration of social dominance and spatial structure. Anim. Behav. 59 (5), 1035–1048.

Hemelrijk, C., 2002. Despotic societies, sexual attraction and the emergence of male 'tolerance': an agent-based model. Behaviour 139 (6), 729–747.

Henrich, J., Boyd, R., Bowles, S., Camerer, C., Fehr, E., Gintis, H., McElreath, R., 2001. In search of homo economicus: behavioral experiments in 15 small-scale societies. Am. Econ. Rev. 91 (2), 73–78.

Herrmann, B., Thöni, C., Gächter, S., 2008. Antisocial punishment across societies. Science 319 (5868), 1362.

Huston, T., Ruggiero, M., Conner, R., Geis, G., 1981. Bystander intervention into crime: a study based on naturally-occurring episodes. Soc. Psychol. Q. 44 (1), 14–23.

Hölldobler, B., Wilson, E., 1990. The Ants. Belknap Press, Cambridge, MA.

Ifti, M., Killingback, T., Doebeli, M., 2004. Effects of neighbourhood size and connectivity on the spatial continuous prisoner's dilemma. J. Theor. Biol. 231 (1), 97–106.

Janssen, M., Bushman, C., 2008. Evolution of cooperation and altruistic punishment when retaliation is possible. J. Theor. Biol. 254 (3), 541–545.

Kagel, J., Roth, A., Hey, J., 1995. The Handbook of Experimental Economics. Princeton University Press, Princeton, NJ.

Killingback, T., Doebeli, M., 1996. Spatial evolutionary game theory: hawks and doves revisited. Proc.: Biol. Sci. 263 (1374), 1135–1144.

Latane, B., Darley, J., 1968. Group inhibition of bystander intervention in emergencies. J. Personal. Soc. Psychol. 10, 215.

Lehmann, L., Keller, L., West, S., Roze, D., 2007. Group selection and kin selection: two concepts but one process. Proc. Natl. Acad. Sci. 104 (16), 6736–6739.

Lewis, D., 1969. Convention: A Philosophical Study. Harvard University Press, Cambridge, MA.

Lindgren, K., Nordahl, M., 1994. Evolutionary dynamics of spatial games. Physica D: nonlinear Phenom. 75 (1–3), 292–309.

Monnin, T., Ratnieks, F., 2001. Policing in queenless ponerine ants. Behav. Ecol. Sociobiol. 50 (2), 97–108.

Mulder, R., Langmore, N., 1993. Dominant males punish helpers for temporary defection in superb fairy-wrens. Anim. Behav. 45, 830–833.

Nakamaru, M., Iwasa, Y., 2006. The coevolution of altruism and punishment: role of the selfish punisher. J. Theor. Biol. 240 (3), 475–488.

Nikiforakis, N., Normann, H., 2008. A comparative statics analysis of punishment in public-good experiments. Exper. Econ. 11 (4), 358–369.

Nowak, M., 2006. Five rules for the evolution of cooperation. Science 314 (5805), 1560.

Nowak, M., 2006. Evolutionary Dynamics: Exploring the Equations of Life. Harvard University Press, Cambridge, MA.

Nowak, M., May, R., 1993. The spatial dilemmas of evolution. Int. J. Bifur. Chaos 3, 35.

Nowak, M., Sigmund, K., 2004. Evolutionary dynamics of biological games. Sci. Signal. 303 (5659), 793–799.

Nowak, M., Tarnita, C., Wilson, E., 2010. The evolution of eusociality. Nature 466 (7310), 1057–1062.

Orbell, J., Dawes, R., 1993. Social welfare, cooperators' advantage, and the option of not playing the game. Am. Sociol. Rev. 58 (6), 787–800.

Oster, G., Wilson, E., 1979. Caste and Ecology in the Social Insects. Princeton University Press, Princeton, NJ.

Ostrom, E., 2000. Collective action and the evolution of social norms. J. Econ. Perspect., 137–158.

Pantin, H., Carver, C., 1982. Induced competence and the bystander effect. J. Appl. Soc. Psychol. 12 (2), 100–111.

Rand, D., Armao, J.J., Nakamaru, M., Ohtsuki, H., 2010. Anti-social punishment can prevent the co-evolution of punishment and cooperation. J. Theor. Biol. 265 (4), 624–632.

Reuben, E., Riedl, A., 2009. Enforcement of contribution norms in public good games with heterogeneous populations. Res. Memo. 29, 1–25.

Semmann, D., Krambeck, H., Milinski, M., 2003. Volunteering leads to rock–paper–scissors dynamics in a public goods game. Nature 425 (6956), 390–393.

Sigmund, K., 2010. The Calculus of Selfishness. Princeton University Press, Princeton, NJ.

Sigmund, K., Hauert, C., Nowak, M., 2001. Reward and punishment. Proc. Natl. Acad. Sci. 98 (19), 10757–10762.

Sigmund, K., DeSilva, H., Traulsen, A., Hauert, C., 2010. Social learning promotes institutions for governing the commons. Nature 466 (7308), 861–863.

Taylor, P., Jonker, L., 1978. Evolutionary stable strategies and game dynamics. Math. Biosci. 40 (1-2), 145–156.

Tomasello, M., 2009. Why We Cooperate. MIT Press, Cambridge, MA.

Trivers, R., Hare, H., 1976. Haplodiploidy and the evolution of the social insect. Science 191 (4224), 249–263.

Warneken, F., Tomasello, M., 2008. Extrinsic rewards undermine altruistic tendencies in 20-month-olds. Dev. Psychol. 44 (6), 1785–1788.

Wenseleers, T., Helanterä, H., Hart, A., Ratnieks, F., 2004. Worker reproduction and policing in insect societies: an ESS analysis. J. Theor. Biol. 17 (5), 1035–1047.

Wenseleers, T., Tofilski, A., Ratnieks, F., 2005. Queen and worker policing in the tree wasp Dolichovespula sylvestris. Behav. Ecol. Sociobiol. 58 (1), 80–86.

West, S., Pen, I., Griffin, A., 2002. Cooperation and competition between relatives. Science 296 (5565), 72–75.

West, S., Griffin, A., Gardner, A., 2007. Social semantics: altruism, cooperation, mutualism, strong reciprocity and group selection. J. Evol. Biol. 20 (2), 415–432.

Wilkinson, G., 1984. Reciprocal food sharing in the vampire bat. Nature 308 (5955), 181–184.

Wilson, E., 1980. Caste and division of labor in leaf-cutter ants (Hymenoptera: Formicidae: Atta). Behav. Ecol. Sociobiol. 7 (2), 143–156.

Wu, J., Zhang, B., Zhou, Z., He, Q., Zheng, X., Cressman, R., Tao, Y., 2009. Costly punishment does not always increase cooperation. Proc. Natl. Acad. Sci. 106 (41), 17448–17451.