

Evolution of Preferences in a Simple 'Game of Life'

Anders Poulsen and Odile Poulsen*

Department of Economics
The Aarhus School of Business

Prismet
Silkeborgvej 2
DK-8000 Aarhus C
Denmark

July 16, 2003

Abstract

The conventional assumption in economics, that individuals have 'materialistic' preferences, has been questioned by experimental evidence. In this paper we endogenize preferences when players are engaged in several different strategic situations. We show that as long as each possible strategic situation is encountered by individuals with a strictly positive probability, there is a unique asymptotically stable preference distribution where reciprocal, altruist and materialist preferences co-exist.

Keywords: Social preferences; preference evolution; reciprocity; altruism; materialism; Prisoner's Dilemma game; simultaneous/sequential moves; multiple games; evolutionary stability.

JEL Classification: B41; C70; C72; D74; Z13.

1 Introduction

The conventional assumption in economics, that individuals are solely motivated by maximization of material returns, has been questioned by carefully controlled lab experiments. Players often behave in a cooperative, or fair, way that does not seem compatible with a self-interested preference. There is, for example, often co-operation in a (sequence of) one-shot Prisoner's Dilemma (PD) games (see e.g. Dawes and

*We thank Stefan Napel for helpful comments. Odile Poulsen gratefully acknowledges financial support from the Danish Social Science Research Council (SSF Grant no. 212.2269.01). This paper is an improved version of Aarhus School of Business Department of Economics Working Paper 02-18, which has circulated under the name "Evolution of Materialism, Reciprocity and Altruism in an Environment of Prisoner's Dilemma Games".

Thaler (1988) and Cooper et. al. (1996)). The same is true for situations where players contribute to a public good (see Ledyard (1995)). In bargaining experiments with games like the ultimatum game (see Güth et. al. (1982)) fairness and reciprocity norms seem to affect the outcomes (see Roth (1995) for an overview). Moreover, there are plenty of examples from everyday life: People give to charity and vote; they return lost wallets; they tip the waiter even when being sure never to visit the same restaurant again, and so on. On the negative side, a customer may refuse to buy a good from a seller in order to punish the seller for charging an 'unfair' high price; an employee who has been sacked may, on her last day at work, engage in acts of sabotage, in order to punish the employer (see Sobel (2001)).

In all these situations people behave in a way that does not maximize their monetary returns, even when they do not have to worry about punishment or future interaction. One possibility is that individual blindly follow social norms and that these completely overwhelm any individual preference. However, clearly not all people are so obedient; and in many of the examples listed above those who decide to violate the norm will make more money than those who follow the social norm. Thus the behavior cannot be explained only by reference to social norms. Another possibility, which we pursue here, is that the observed behavior is due to people's internalized personal norms, embodied in their *preferences*: People differ in their preferences and this is why they behave differently, notwithstanding any social norms. Indeed, experimental research has documented the existence of a significant proportion of players with 'social preferences'. Following Fehr and Fischbacher (2002), we say that a player has social preferences when she cares not just about her own money payoffs, but also cares about the entire distribution of money payoffs between her and other reference agents, and/or cares about how this distribution was brought about. We refer the reader to e.g. Fehr and Schmidt (2001) for a survey of the experimental findings and for references.

There is considerable empirical evidence indicating that a significant proportion of people have social preferences and that another significant proportion have materialistic preferences (see e.g. Fehr and Gächter (1998)). We would like to explain why there apparently is such a distribution of preferences in the population, with some people having social preferences, and others materialistic preferences. In this paper we therefore endogenize the proportions of individuals having different preferences. To do this, we let different preferences compete against each other, in order to determine what preference(s) will emerge as the 'winner' of the economic struggle for survival. This methodology is called the 'indirect evolutionary approach' (see e.g. Güth and Yaari (1992) and Güth (1995)): Players act rationally given their preferences, but those preferences may change over time, as the result of a socioeconomic or cultural learning and imitation process. The key assumption is that preferences who give their human carriers higher-than-average *money* payoffs tend to be adopted by more players over time. An interpretation is that, in order to survive in the economic system, one needs to perform *materially* well. However, this assumption, that 'only money matters', does not a priori bias the analysis towards the survival of materialistic preferences: *Any* sort of preference that leads players to a materially superior behavior will prosper. Thus, if players with, say, reciprocal, preferences earn more money than the average, and those with materialistic preferences less, the proportion of the former individuals will increase at the expense of the latter.

After having set the stage for competition between preferences, the next step is to outline the strategic environment in which the competition will take place. Needless to say, this is crucial since the monetary returns that a preference confer on an individual depends on the features of the strategic environment within which the preference is operating. A preference may generate behavior that is very successful in one environment, but may give rise to inferior behavior in another situation. We will refer to the multitude of strategic situations that individual encounter as the *Game of Life*.¹ Our goal is to see how the nature of the endogenously determined preferences depend on the features of this Game of Life. Most existing models of preference evolution have simplified the analysis by assuming that the Game of Life consisted of a single game only, such as a duopoly game (Bester and Güth (1998), Bolle (2000) and Possajennikov (2000)), a Prisoner's Dilemma (Fershtman and Weiss (1998), Engelmann (2001), Guttman (2000), Possajennikov (2002)), a Trust Game (Güth and Yaari (1992) and Güth (1995), Güth and Kliemt (1998)), an Ultimatum Game (Huck and Oechssler (1999)), A 'Hawk-Dove' (or 'Chicken') game (Possajennikov (2002), Poulsen (2001)) and several other games. There are, to our knowledge, very few papers allowing for more than a single situation; in Güth and Napel (2002) the Game of Life consists of an Ultimatum Game and a Dictator game. Samuelson (2001) also considers a variety of game situations, but uses it to explore the implications of bounded rationality.

The contribution of this paper is to study the role of the following aspects of the Game of Life: How do people interact - simultaneously or sequentially? In the Game of Life people are involved in both kinds of situations, and hence they will evolve preferences that perform 'well', as an average over these game situations. Precisely, we assume that individuals encounter either a simultaneous one-shot Prisoner's Dilemma or a sequential Prisoner's Dilemma; which one is determined exogenously with fixed probabilities. The distinction between simultaneous or sequential interaction is relevant because different preferences experience different advantages and disadvantages in either context. For example, some preferences dictate choices that are conditional on the opponent's moves. This is the case for important preferences such as reciprocity (see Rabin (1993)) and inequality aversion (Fehr and Schmidt (1999)). When interaction is simultaneous individuals equipped with such preferences often face multiple equilibria, with a 'good' (payoff-dominant) equilibrium and an inferior equilibrium. Typically, it is assumed that the individuals manage to select the payoff-dominant equilibrium with probability one (see for example Guttman (2000) and Engelmann (2001)). As we shall see, however, whenever there is the tiniest risk that two individuals with reciprocal preferences fail to co-ordinate on the payoff-dominant equilibrium, then the reciprocal preference is at a disadvantage compared with simpler preferences who simply dictate an unconditional choice. Another paper that allows for less than perfect co-ordination is Fershtman and Weiss (1998). Conversely, the advantage of evolving simple preferences is that one, when moving simultaneously with similar individuals, avoids any co-ordination problems. On the other hand, individuals with preferences that specify simple unconditional behavior are at risk of being exploited by other individuals, and/or may perform badly against each other.

We show that, whenever there is some, no matter how small, difficulty in overcoming the co-ordination problems mentioned earlier, then preference evolution, seeking to balance the relative advantages of the various preferences in the simultaneous and se-

¹We have borrowed this term from Güth and Napel (2002).

quential interactions involved the Game of Life, leads to a unique asymptotically stable preference distribution. Three preferences co-exist in the endogenous preference distribution: A 'Materialist', a 'Reciprocator' and an 'Altruist' preference. The Materialist always defects, the Altruist always co-operates and the Reciprocator co-operates (defects) if the opponent co-operates (defects). Thus there will be both co-operation and defection in the population. We obtain asymptotic stability for any non-degenerate probability distribution over the simultaneous and the sequential game. The three preferences co-exist because they each have certain advantages in the Game of Life: The Reciprocator preference survives because it avoids exploitation and performs relatively well against itself; still, as long as there are some co-ordination problems coming from simultaneous interaction, it cannot alone dominate the population; the Altruist preference survives because it induces Reciprocator to treat it well and because Altruists perform well against each other; and, finally, the Materialist preference survives because it performs very well against Altruist individuals. Moreover, allowing for multiple games actually strengthens the predictions of the model, since, as we also show, the prediction for each of the component games of the Game of Life is not unique.

We also analyze the case where two reciprocal individuals choose the co-operative equilibrium with probability one in the simultaneous game. In this case we get co-existence between the reciprocal and altruist preference only. However, instead of studying this cases from the outset, we argue that it is more plausible to *first* compute the endogenous preference distribution when Reciprocators can establish *some* co-operation and *then* consider the limiting case where reciprocal individuals become perfectly able to select the co-operative equilibrium. Indeed, whereas the limiting preference distribution is unique, when Reciprocators perfectly co-ordinate on the co-operative the predicted preference distributions are not unique.

Our results are based on an assumption, common in the literature, that an individual learns, prior to interaction, what the opponent's preferences are. We also investigate what effect it has on the Game of Life that players receive no information at all about fellow individuals' preferences. We show that in this case social preferences have no impact on behavior: In any evolutionarily stable outcome, all players defect.² In communities with anonymous interaction, where players obtain little information about fellow individuals' preferences, reciprocity and altruistic behavior should not be expected to emerge. If, on the other hand, it is possible for players to acquire information about opponents prior to interaction, then reciprocity and/or altruism is a real possibility, and sometimes the only possibility. These results should, in our opinion, be interpreted as an indication, in an evolutionary context, that the standard assumption that players *always and only* have materialistic motivations, is ill-founded and inappropriate. Such a statement needs to be qualified: Certain environments are more conducive for social preferences than others. The features that we focus on are the way people interact (simultaneous versus sequential play) and the amount of information that is available about other people's preferences. But, how is it that a player can deduce whether another person is, say, a reciprocal or materialist individual? One possibility is that individuals have access to information about an opponent's previous behavior, from encounters with other people. Given this information, people form a belief about the opponent's type (preference). See e.g. Kandori (1992) for such a

²More general results along those lines are given in Dekel et. al. (1998) and Ok and Vega-Redondo (2001).

model. Yet another possibility is that individuals base their evaluations of other individuals' preferences using indicators such as income, skin color, area of residence, and so on. We leave for future research the task to incorporate such, and other, realistic features into models of preference evolution.

What is the relevance of our results for the experimental findings, mentioned earlier? In most experiments players do not have information about fellow players' preferences. Interaction is deliberately kept anonymous. Our results, emphasizing the role of information about preferences for reciprocity and altruism to survive, can therefore not *directly* explain the experimental findings. What happens in the lab is presumably that subjects' believe that there are sufficiently many reciprocal and altruist people out there. Given these (correct) beliefs subjects with reciprocal/altruistic preferences optimally cooperate. The question is, therefore, once more: How is it that some (often a significant proportion of) subjects have these beliefs and preferences? It is this question that our model supplies some answers to: The beliefs and preferences have been shaped in the outside 'Game of Life' and are consequently used in the experimental lab, too. The Game of Life is such that players with social preferences do survive, side by side with materialistic players. In our model not all players could be materialists in the Game of Life, for reciprocal players could invade. Similarly, not everybody could be reciprocal, or altruist, since players with other preferences would outperform them. The result is co-existence between players with different preferences. This is why some players, when seated in the lab in an anonymous setting, optimally co-operate, while others do not. We do not, of course, claim that our simple Game of Life, modeled as a mix of simultaneous and sequential Prisoner's Dilemma games, is an adequate representation of the real Game of Life. However, we still believe it gives an insight into what one might expect from a richer model.

There are, in addition to the papers already mentioned earlier, several other models of preference evolution, studying different games and using somewhat different modeling techniques. We refer the reader to Dekel, Ely and Yilankaya (1998), Ely and Yilankaya (2001), Ok and Vega-Redondo (2001), Possajennikov (2002), Sandholm (2001) and Sethi and Somanathan (2001). In Dekel et. al. (1998), for example, preference evolution operates on a symmetric normal form game. Their analysis is very general and based on the static evolutionary solution concept of a Neutrally Stable Strategy, or NSS (Maynard Smith (1982)). They obtain results on the relationship between efficiency and the evolutionarily stable preference distributions. In particular, they show that stability implies efficiency (when preferences are perfectly observable; they also analyze the cases of imperfect and no information about preferences); a similar result is obtained in Possajennikov (2002). This is essentially due to a 'secret handshake' argument; we refer the reader to Robson (1990). This 'stability implies efficiency' result does not hold in our model. The reasons are several. First, in addition to studying general symmetric normal form games, Dekel et. al. allow for a much larger set of preference types than we do. For example, one feasible preference type in their model (and in Possajennikov (2002)) is the type that is indifferent between all her feasible actions no matter what the opponent does; this type is not feasible in our model. It is the inclusion of such preference types that impose restrictions on what preferences can be evolutionarily stable in their model. We do not admit such preferences, since our research purpose is very different from theirs. Second, it is well-known that using the static NSS solution concept is coarser than directly studying the stable

outcomes for a given dynamic, such as the Replicator Dynamic; a population can be stable for this dynamic but not be an NSS (see e.g. Weibull (1995)). One of our stable preference distributions display exactly this feature (cf. Proposition 3 below).³

The rest of the paper is organized as follows. In Section 2 we set up the basic model, specifying the component games in the Game of Life, preferences, information and the evolutionary selection of preferences. In Section 3 we first study the component games separately. Then, in Section 4 we study the case where both simultaneous and sequential interaction can occur with arbitrary positive probabilities. In Section 5 we allow for more feasible preferences. We study the case where players receive no information about other players' preferences in Section 6. The robustness of our results and certain assumptions is investigated in Section 7. We conclude in Section 8. All proofs are in the Appendix.

2 The Model

2.1 The Environment

There is a large population of players. Time is continuous and at each time instant all players are randomly matched in pairs. With probability μ , a pair of individuals play a simultaneous Prisoner's Dilemma (PD) game, with the following *money* payoffs:

	<i>C</i>	<i>D</i>
<i>C</i>	1 <i>b</i>	
<i>D</i>	<i>a</i> 0	

where $a > 1 > 0 > b$, $(1/2)(a + b) < 1$ and '*C*' and '*D*' stand for 'Co-operate' and 'Defect', respectively. With complementary probability, $1 - \mu$, the pair plays a sequential PD game. Here one of the players is randomly chosen to be first-mover. This player then chooses between co-operate and defect. The other player, the second-mover, observes the first-mover's choice and makes a choice himself. The money payoffs from the various outcomes in the sequential game are the same as in the matrix above. The fact that players in the sequential game are randomly allocated to be first-mover or second-mover is meant to reflect a set-up where, when two players meet each other, random factors decide who moves first.

2.2 Preferences

Let (i, j) , where $i, j = C, D$, denote the outcome where a player chooses i and the opponent chooses j . We will consider the (pure) *best replies*, i.e., what a player will choose if the opponent plays C , and what he will choose if the opponent plays D .

In our model, a *preference type* is a pair of pure best replies. We will consider all such pairs:

³This raises the issue whether, in the model of Dekel et. al., an inefficient preference distribution could be stable for the Replicator Dynamic, even though it is not an NSS.

The **Materialist** (M) preference type: Choose D both if the opponent chooses C and if the opponent chooses D . That is, D is strictly dominant. The **Reciprocator** (R) preference type: Play C if the opponent chooses C and play D if the opponent chooses D . The **Altruist** (A) preference type: Choose C both if the opponent chooses C and if the opponent chooses D . That is, C is strictly dominant. The **Paradoxical** (P) preference type: Choose D if the opponent chooses C and choose C if the opponent chooses D .

A player is characterized by one of these four preference types. It is important that a player cannot evolve a separate preference type for each of the situations he is involved in. In principle, a preference type could be a triple (i, j, k) , where $i, j, k = M, R, A, P$, i is the preference type of the individual when he is playing the simultaneous PD and j (k) is the preference type when he is first-mover (second-mover) in the sequential game. However, we do not allow for this, since our purpose is exactly to expose a preference to several game situations.

We will assume, until Section 6, that preferences are common knowledge. For expositional reasons we will initially exclude the Paradoxical preference type from the analysis. Indeed, in Section 5 we show that this preference type under certain conditions will go extinct.

The four preference types are the primitives of our model. Of course, in reality the primitives are the players' orderings of the outcomes. However, many of these preference orderings will induce the same behavior in the Game of Life; hence they may be regarded as equivalent and there is no need to distinguish between them. However, this not true for other orderings. Nevertheless, in Section 7 below we show that our results are, in a certain sense, the same as those one would obtain from considering the entire set of preference orderings instead of the four best reply pairs.

2.3 Endogenizing Preferences

Let x_i , with $i = A, R, M, P$, denote the population proportion of players of type i , where $0 \leq x_i \leq 1$ and $\sum_i x_i = 1$. Then the state of the system is the *preference distribution* $x = (x_A, x_R, x_M)$. Denote by $\pi(i, x)$ the expected *money* payoff to a type i player and let $\pi(x, x)$ denote the average expected money payoff at the preference distribution x . Then the evolution of the proportion of players of type i is given by

$$\dot{x}_i = x_i[\pi(i, x) - \pi(x, x)].$$

This is the Replicator Dynamic (Taylor and Jonker (1978)). It says that the growth rate of players with preference $i = A, R, M, P$ is positive if these players earn above-average *money* payoff. We wish to describe the dynamic of preference evolution and to find those preference distributions that are (asymptotically) stable for this dynamic.

3 A Benchmark: Either Simultaneous or Sequential Interaction

In this section we analyze, as a benchmark, the simultaneous and sequential game on their own. As already mentioned in Section 2.2, we will not include the Paradoxical preference type until Section 5 below.

3.1 Simultaneous Interaction

For the evolutionary analysis we need to compute the money payoffs π_{ij} , where $i, j = A, R, M$, obtained by a player with preference i when she is matched with an opponent of type j . These money payoffs are given in the matrix below⁴:

	A	R	M
A	1	1	b
R	1	π_{RR}	0
M	a	0	0

Table 1: The money payoffs in the evolutionary game under simultaneous interaction. A = Altruist; R = Reciprocator; M = Materialist.

In a meeting between two R -types, there are two possible outcomes, corresponding to the two strict Nash equilibria: (D, D) and (C, C) .⁵ In accordance with the discussion in the Introduction, we make the following assumption, also employed by Fershtman and Weiss (1998):

Assumption 1 *Two Reciprocators select the (C, C) equilibrium with probability λ , and select the (D, D) Nash equilibrium with probability $1 - \lambda$, where $0 < \lambda < 1$.*

Below, in Section 4.1, we also consider the cases $\lambda = 1$ and $\lambda = 0$, in which case the conclusion is very different. Indeed, for reasons that will be made clear below we prefer first to compute the endogenous preference distribution for some fixed $0 < \lambda < 1$ and *then* consider the limiting preference distribution, as $\lambda \rightarrow 1$ or $\lambda \rightarrow 0$, rather than setting $\lambda = 1$ or $\lambda = 0$ from the outset.

With Assumption 1 we have

$$0 < \pi_{RR} = \lambda < 1. \tag{1}$$

We then obtain the following result:

⁴Consider, for example, a meeting between an M -type and an R -type. The M -type always plays D and the R -type consequently plays D , too. Thus the money payoff to each player is the mutual defection payoff, zero: $\pi_{MR} = \pi_{RM} = 0$. Similarly, in an encounter between an A -type and an R -type, the former always plays C and the R -type then responds with C , too. Thus $\pi_{AR} = \pi_{RA} = 1$.

⁵There is also a symmetric and mixed Nash equilibrium, but we ignore it here.

Proposition 1 Consider the evolutionary game based on the simultaneous-move Prisoner's Dilemma game. There is a unique stable equilibrium preference distribution, x^* . This is an interior equilibrium, i.e., $x_i^* > 0$ for all $i = A, R, M$. x^* is a center, i.e., it is surrounded by closed orbits.

$$x^* = (x_A^*, x_R^*, x_M^*) = \left(\frac{-b\pi_{RR}}{D}, \frac{b(1-a)}{D}, \frac{(1-\pi_{RR})(a-1)}{D} \right), \quad (2)$$

where $D = (1-a-b)\pi_{RR} + (a-1)(1-b)$ and $\pi_{RR} = \lambda$.

Proof: Please see the Appendix.

Remark: The proof that the preference distribution x^* is a center follows from applying results in Hofbauer (1981) and Bomze (1983). We may also note that x^* is neither an Evolutionarily Stable Strategy, nor a Neutrally Stable Strategy (see Weibull (1995)). This is because a mutant equipped with the Materialist preference is an alternative best reply to x^* and violates both definitions' stability conditions. Nevertheless x^* is stable.

The crucial features giving the result in Proposition 1 are that two Reciprocators cannot perfectly co-ordinate on the co-operate equilibrium ($\lambda < 1$), which implies that an Altruist invades the all-Reciprocator population; and that two Reciprocators can establish some co-operation ($\lambda > 0$), which makes the all-Materialist population unstable.

3.2 Sequential Interaction

We now consider the sequential game. Consider a player who, in the role of first-mover, faces a reciprocal second-mover. If the first-mover cooperates (defects), the second-mover cooperates (defects), too. Thus we must ask how the first-mover ranks the (C, C) outcome relative to the (D, D) outcome.

Assumption 2 All players, irrespective of their preference type, prefers the (C, C) outcome to the (D, D) outcome.

In Section 7 we show that evolutionary selection on a wider class of preference types will make sure that this assumption is satisfied. We believe this is a plausible assumption, which is entirely in line with our proposed interpretation of materialism, reciprocity and altruism.⁶ We obtain the following matrix for the evolutionary game:

	A	R	M
A	1	1	b
R	1	1	$1/2$
M	a	$1/2$	0

⁶Indeed, it seems difficult to find an otherwise plausible utility function that would make (D, D) preferred to (C, C) .

Table 2: The money payoffs in the evolutionary game under sequential interaction. A = Altruist; R = Reciprocator; M = Materialist.

We see that (i). sequential interaction allows two Reciprocators to overcome their co-ordination problem and (ii). a Reciprocator outperforms a Materialist against a Materialist opponent. Observation (i) holds because when a reciprocal first-mover faces a reciprocal second-mover, Assumption 2 implies that the first-mover chooses C and the second-mover then responds with C . In a meeting between two Reciprocators, each player therefore earns $(1/2)(1) + (1/2)(1) = 1$. The second observation comes from the fact that when the Materialist is first-mover she optimally chooses to *cooperate* rather than to defect (Assumption 2). Thus the presence of reciprocally motivated players induces materialistic players to behave more cooperatively than they otherwise would. Even though the Materialist prefers to defect, *given* any choice by the opponent, when her own choice will determine the opponent's choice, she is led to cooperate. However, when the Reciprocator is first-mover she defects when facing a Materialist second-mover. In this sense, the Materialist is more co-operative than the Reciprocator in sequential interaction. This is a quite general phenomenon: The presence of reciprocal individuals in the population affects the behavior of materialistic individuals; and whenever both types are present, the former (latter) player type becomes less (more) co-operative. See e.g. Fehr and Schmidt (1999) and Fehr and Fischbacher (2002) for a discussion. We obtain the following proposition:

Proposition 2 *Consider the evolutionary game for the Sequential Prisoner's Dilemma game. Any preference distribution composed exclusively of Reciprocators and Altruists and where the proportion of Reciprocators is sufficiently large, $x_R > 2(a - 1)/(2a - 1)$, is stable.*

Proof: Please see the Appendix.

There are no Materialists in any stable population. The Materialist's expected payoff is higher in the sequential game than in the simultaneous game. However, the Reciprocator's is even higher: Not only does sequential interaction allow the Reciprocators to overcome their internal' equilibrium selection problem, but they outperform the Materialists against other Materialists. All this implies that a Reciprocator preference can be sustained, and this, in turn, allows for some altruism to survive, too.

4 Both Simultaneous and Sequential Interaction

In this section we assume that when two players are matched, they engage in the simultaneous game with probability μ and with remaining probability $1 - \mu$ interaction is sequential. We will denote the resulting game as the ' μ -game'. We impose only $0 < \mu < 1$. A player's evolutionary performance, i.e., his monetary earnings, is now a weighted average of his performance in the simultaneous game and his performance in the sequential one. With respect to the simultaneous game, we maintain Assumption 1 (we actually only need $\lambda < 1$; see the discussion below). We then obtain the following matrix for the evolutionary game:

	A	R	M
A	1	1	b
R	1	$\mu\pi_{RR} + 1 - \mu$	$(1 - \mu)(1/2)$
M	a	$(1 - \mu)(1/2)$	0

Table 3: The money payoffs in the μ -game where the simultaneous (sequential) game is played with probability μ ($1 - \mu$).

What are the crucial features that the μ -game inherits from the simultaneous and the sequential game? Let us recall that (i). in the simultaneous (sequential) game preference type A is the unique (an alternative) best reply to preference type R , (ii). in both the simultaneous and the sequential game the M preference type is the unique best reply to preference type A and (iii). in the sequential (simultaneous) game preference type R is the unique (alternative) best reply to preference type M . The implication is that for any $\mu \in (0, 1)$ we get a cyclical best reply structure: A is a unique best reply to R , M is a unique best reply to A and R is a unique best reply to M . We did not have this cyclical relationship in the simultaneous or the sequential games when analyzed on their own. This gives rise to a dynamic and stability property that is qualitatively from those of the component games:

Proposition 3 *Consider the evolutionary μ -game, where two individuals play the simultaneous Prisoner's Dilemma game with probability μ and the sequential game with probability $1 - \mu$, where $0 < \mu < 1$. There is, for each $\mu \in (0, 1)$ and $\lambda \in [0, 1)$, a unique interior equilibrium, $y^* = (y_A^*, y_R^*, y_M^*)$, and it is globally asymptotically stable.*

$$y_A^* = \frac{\mu^2 + 2b(2\pi_{RR} - 1)\mu + 2b - 1}{E}, \quad (3)$$

$$y_R^* = \frac{2(a - 1)(2b - 1 + \mu)}{E}, \quad (4)$$

$$y_M^* = \frac{4\mu(1 - a)(1 - \pi_{RR})}{E}, \quad (5)$$

where $E = 4\mu\pi_{RR}(a + b - 1) + (2b - 1 - \mu)(2a - 1 - \mu)$ and $\pi_{RR} = \lambda$.

Proof: Please see the Appendix.

The dynamic, with $a = 2$, $b = -1$ and $\lambda = 1/2$, and with each game being equally likely to be played, i.e., $\mu = 1/2$, is shown below.

Evolutionary selection for the μ -game gives a unique prediction of population behavior for each value of $\mu \in (0, 1)$: Any initial preference distribution, where all three types are present, will over time evolve to the preference profile y^* . In the simultaneous game, on the other hand, we got perpetual cycles, and in the sequential game there was a entire set of stable populations. Thus, by allowing for multiple games being played, we obtain a sharper prediction than in either the simultaneous or the sequential game.

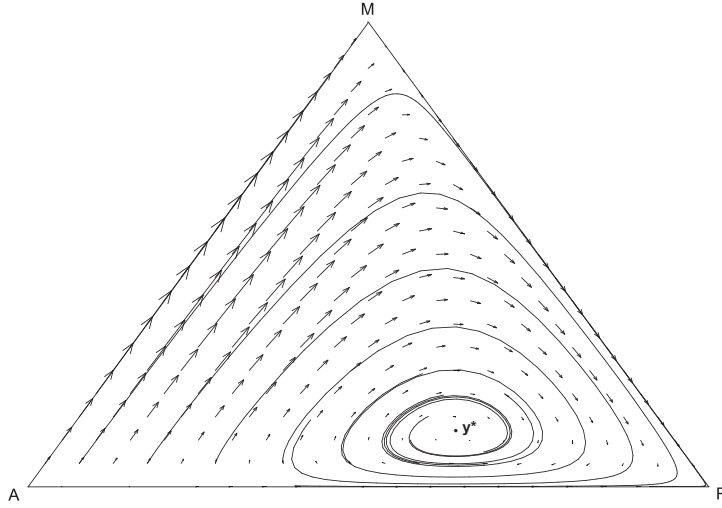


Figure 1: Phase diagram. Parameter values: $a = 2$, $b = -1$, $\pi_{RR} = 1/2$ and $\mu = 1/2$. Equilibrium proportions: $y_A^* = 11/35$, $y_R^* = 20/35$ and $y_M^* = 4/35$.

4.1 Imposing $\lambda = 1$ from the Outset

Let us now directly set $\lambda = 1$, i.e., two Reciprocators select the payoff-dominant (C, C) equilibrium with probability one. In this case there will not be a unique stable preference distribution, as shown by the following result.⁷

Proposition 4 *Suppose two Reciprocators play the (C, C) Nash equilibrium with probability one in the simultaneous game ($\lambda = 1$). Then there is no interior stable preference distribution. Instead altruism survives with reciprocity for any $\mu \in (0, 1)$, as long as there are sufficiently few Altruists. Precisely, any population with Reciprocators and Altruists only, and where the proportion of Reciprocators is no smaller than $\max\{(a - 1)/a, (a - 2)/(a - 1 - b)\}$, is a Neutrally Stable Strategy (NSS) and hence stable for the Replicator Dynamic for any $\mu \in (0, 1)$.*⁸

Proof: Please see the Appendix.

This result is similar to Proposition 2. The reason is that when Reciprocators can co-ordinate flawlessly on the co-operative equilibrium in the simultaneous game, the simultaneous and sequential game have the same qualitative evolutionary stability properties, and the Game of Life collapses to a single game. The reason is that the all-Reciprocator preference distribution can no longer be invaded by Altruists; the latter can, however, still enter the population. As long as there are sufficiently few Altruists, Materialists do not earn enough to invade the population.

4.2 Some Limiting Results

From Proposition 3, we obtain the following results:

⁷Similar results are given in Poulsen (2003).

⁸We refer the reader to Weibull (1995) for the NSS concept, due to Maynard-Smith (1982).

$$\lim_{\mu \rightarrow 1} y^* = x^* \quad (6)$$

$$\lim_{\mu \rightarrow 0} y^* = \left(\frac{1}{2a-1}, \frac{2(a-1)}{2a-1}, 0 \right). \quad (7)$$

The first result is not surprising: In the limit, as players are only engaged in the simultaneous game, the preference distribution collapses to the one from Section 3.1. The second one is perhaps more interesting: In the limit, as only the sequential game is played, the preference distribution approaches a unique distribution where individuals are either Altruists or Reciprocators. This is a sharper result than what Proposition 2 could deliver.⁹

Considering the effects of varying λ , we obtain

$$\lim_{\lambda \rightarrow 1} y^* = (p, q, 0), \quad (8)$$

where $p > 0$ and $q = 1 - p > 0$ and are obtained from Proposition 3. Moreover,

$$\lim_{\lambda \rightarrow 0} y^* = \left(\frac{\mu^2 - 2b\mu + 2b - 1}{F}, \frac{2(a-1)(2b-1+\mu)}{F}, \frac{4\mu(1-a)}{F} \right), \quad (9)$$

where $F = (2b-1-\mu)(2a-1-\mu)$. (9) shows that even if, in the simultaneous game, the Reciprocators were to become unable to establish any co-operation between them at all, we would still, for $\mu < 1$, obtain an interior preference distribution. In other words, in the μ -game what matters is not whether $\lambda = 0$ or $\lambda > 0$, but whether $\lambda < 1$ or $\lambda = 1$. Whereas the limiting preference distribution (8) is unique, our prediction from setting $\lambda = 1$ from the outset were indeterminate (cf. Proposition 4). Thus, the method of first computing the equilibrium preference distribution for some λ with $0 < \lambda < 1$ and *then* considering the limiting preference distribution, as $\lambda \rightarrow 1$, yields sharper predictions than the one obtained from setting $\lambda = 1$ from the outset.

5 Including the 'Paradoxical' Preference Type

Until now we ignored the Paradoxical preference type: Play D if the opponent plays C and play C if the opponent plays D . Even though this may look like a rather strange preference, we should not be too quick in excluding it.

We proceed straight to the μ -game. The matrix below contains all the money payoffs for the evolutionary game.

⁹Indeed, the limiting preference distribution corresponds exactly to the critical distribution mentioned in Proposition 2.

	A	R	M	P
A	1	1	b	b
R	1	$\mu\pi_{RR} + 1 - \mu$	$(1/2)(1 - \mu)$	$\mu\pi_{RP} + (1/2)(1 - \mu)(a + 1)$
M	a	$(1/2)(1 - \mu)$	0	a
P	a	$\mu\pi_{PR} + (1/2)(1 - \mu)(1 + b)$	b	$\mu\pi_{PP} + (1/2)(1 - \mu)(a + b)$

Table 4: The money payoffs in the evolutionary μ -game with four preference types and where the simultaneous (sequential) Prisoner's Dilemma game is played with probability μ ($1 - \mu$); A = Altruist; R = Reciprocator; M = Materialist; P = Paradoxical preference type.

Let us briefly explain how we have arrived at the money payoffs in the fourth row and column. The first payoff, a , is computed as follows. In the simultaneous game the A type plays C and so the P -type responds with D . Thus the P -type gets a . In the sequential game the P -type chooses D when first-mover and the A -type responds with C . When the A -type is first-mover, she effectively chooses between (C, D) , giving her money payoff a , and (D, C) , giving money payoff b . It is inconsequential for our result below what the Altruist prefers in this situation, and so we assume that the Altruist chooses C , thereby establishing the (C, D) outcome.

When the P -type meets an R -type, there is a unique symmetric and mixed Nash equilibrium in the simultaneous game, giving money payoff π_{PR} to the P type. In the sequential game a reciprocal first-mover in effect chooses between the outcomes (D, C) and (C, D) . We make the following assumption:

Assumption 3 *The R type prefers (D, C) to (C, D) .*

Intuitively, this will hold whenever the Reciprocator type cares sufficiently much about his own money payoff and/or does not like to get less than the opponent, relative to being altruistic. In Section 7 below, where we check the robustness of our analysis, we show that the assumption will be satisfied when evolution has as larger menu of preferences available.

It follows that in the sequential game the reciprocal first-mover chooses D and the paradoxical second-mover responds with C . Finally, when the P -type is first-mover, we get the (C, C) outcome. The P -type therefore gets money payoff $\mu\pi_{PR} + (1 - \mu)[(1/2)b + 1/2]$. The R -type gets $\mu\pi_{RP} + (1 - \mu)[(1/2)a + 1/2]$. When the opponent is an M -type, on the other hand, the outcome is that the P -type chooses C and the M -type chooses D , both in the simultaneous and the sequential game. Finally, suppose two P -types meet. In the simultaneous game there is a unique symmetric mixed Nash equilibrium¹⁰, giving money payoff π_{PP} . In the sequential game a P -type as first-mover gets a and as second-mover he gets b . Thus the overall monetary payoff to a P -type against another P -type in the μ -game is $\mu\pi_{PP} + (1 - \mu)[(1/2)a + (1/2)b]$.

Looking at the money payoffs in the matrix reveals that the M type performs strictly better, or at least as well, as the P type against the A , M and P types. When is the P also outperformed against the R type? First, M performs better in the

¹⁰There are also two asymmetric pure Nash equilibria, but we ignore them here.

sequential game, since the P type's preferences lead her to being exploited by an R type in the sequential game. The M -type, on the other hand, realizes the better (D, D) outcome when matched with an R -type and the former is second-mover. However, the P type performs strictly better than the M type against an R type in the simultaneous contest. Thus, for M to outperform P against R it is necessary that interaction is sufficiently sequential, or that P 's advantage in the simultaneous game is sufficiently small:

This gives us the following proposition:

Proposition 5 *Suppose it is sufficiently likely that interaction is sequential:*

$$\mu < \frac{-b}{2a - b}.$$

Then (a). The M -type weakly dominates the P -type. (b). When the initial preference distribution contains all four preference types, the population proportion of type P players approaches zero as time approaches infinity.

Proof: Please see the Appendix.

The condition in the proposition is likely to hold whenever the P type performs badly against the R type in the simultaneous-move setting (such that π_{PR} is small). This, in turn, will be the case whenever the R -type's preferences lead her to behave 'aggressively' against the P type, i.e., to be very likely to play defect.

Part (b) implies that we can effectively ignore the P preference type from the analysis: The preference distribution will eventually 'land' on the face of the simplex spanned by the A , R and M vertices, and from then on the dynamic will be as when only these three strategies were available from the beginning. In this case all our results from Sections 3 and 4 hold. If, on the other hand, interaction is sufficiently likely to be simultaneous, then the P type performs better than M in both the sequential and the simultaneous game. Then it is possible that the Paradoxical preference survives with some of the other preferences. For a precise characterization a more thorough study of dynamic stability is required.

6 The Case of No Information about Preferences

In the previous analysis subjective payoffs were common knowledge. It was as if players had their preference types written on their foreheads. We conjecture that our results continue to hold as long as information is *sufficiently* accurate, or as long as it is not too costly to acquire such information. However, let us now consider the polar opposite to perfect information: A player, when having to decide between cooperating and defecting, receives no information about the opponent's preferences; all interaction is completely anonymous. The only thing a player knows is the *aggregate* distribution of the different preference types in the population.

Anonymous interaction means that a player's preferences can no longer affect an opponent's choice. All players face the same distribution of C and D choices. But

that implies that the unambiguously best thing to do in terms of money is to defect. The following result holds no matter whether interaction is simultaneous, sequential or combined:

Proposition 6 *Consider the simultaneous, the sequential or the μ - game when players do not know their opponents' preferences. Then: In any stable preference distribution there are no Altruists and all players defect.*

Note that the proposition does not say that all players are materialistic in any stable preference distribution; some may be reciprocal, but none altruistic. However, there are so many Materialists that the Reciprocators defect, too. Thus materialistic and reciprocal players are indistinguishable from each other. It is the lack of a means of communication that prevents the Reciprocators from 'breaking out' and establishing the cooperation between themselves that would give them an evolutionary advantage over the materialists. In the terminology of Robson (1990), the reciprocal players cannot give each other a 'secret handshake'.

This proposition underlines the importance of communication for social preferences to have survival value. Players must be able to signal, or communicate, what value system they have, in order to establish cooperative outcomes and in order to avoid co-operating with 'bad' players. In a *completely* anonymous world, social preferences do not make a behavioral difference relative to the case where all players have the usual, materialistic, preferences. Similar results, in more general games, are reported in Dekel et. al. (1998) and Ok and Vega-Redondo (2001). The assumption of no information seems, however, to be as unrealistic as assuming perfect information. In the real world, individuals form opinions about fellow individuals' preferences based on information about observables such as income and skin color. The future challenge, we believe, is to explore this 'intermediate' are, between perfect and no information.

7 Best Replies versus Preference Orderings

Above we defined a 'preference type' as a pair of pure best replies. As remarked in Section 2.2 an alternative approach is to use a set of preference orderings and then obtain the best reply behavior from those orderings. Would our results change if we followed this approach instead? Moreover, are our Assumptions 2 and 3 for the sequential game restrictive? In this section we show that our four preference types can be considered as a reduced form, or approximation to, a more rich situation, and that, additionally, the two above-mentioned assumptions will be satisfied.

There are $4! = 24$ strict preference orderings over the four outcomes; of those six give best replies of the Reciprocator type¹¹, six give those of the Altruist type, and so on. A 'preference type' will now refer to such an ordering. We will say that a preference type is an i - type when the best replies are those of the i -type, where $i = A, R, M, P$. In principle we could study the evolutionary game, and, based on

¹¹The six orderings are: $(C, C) \succ (D, C) \succ (D, D) \succ (C, D)$, $(C, C) \succ (D, D) \succ (D, C) \succ (C, D)$, $(C, C) \succ (D, D) \succ (C, D) \succ (D, C)$, $(D, D) \succ (C, D) \succ (C, C) \succ (D, C)$, $(D, D) \succ (C, C) \succ (C, D) \succ (D, C)$ and $(D, D) \succ (C, C) \succ (D, C) \succ (C, D)$.

the resulting 24×24 payoff matrices (one for the simultaneous and the other for the sequential game), find the stable preference distributions. However, some preference types are behaviorally equivalent and others are weakly dominated. We will use this to narrow down the set of preference types to four, each giving one of the four best reply combinations we have been studying in the previous sections.

We first need to compute the monetary payoffs obtained by the various preference types. Whenever any two Reciprocator types meet each other in the simultaneous game, the same equilibrium selection problem occurs as studied earlier. Moreover, whenever any Reciprocator types meets any Paradoxical type in the simultaneous game there is again a unique mixed and symmetric Nash equilibrium, as is the case whenever any two Paradoxical types meet. We will make the following simplifying assumption, based on the principle of insufficient reason.

Assumption 4 *In the simultaneous game, (i). All R-types have the same expected money earnings against any given R or P type. (ii). All P-types have the same expected money earnings against any given R or P type.*

We then have the following result:

Lemma 1 *Any preference type with $(D, D) \succ (C, C)$ is weakly dominated.*

The proof is quite simple: Consider an ordering with $(D, D) \succ (C, C)$ and compare it with another ordering with the same best replies, but with $(C, C) \succ (D, D)$. From Assumption 4 these two orderings earn the same expected payoffs in the simultaneous game (since only the best replies matter). In the sequential game the same will be true, except when being first-mover and facing a Reciprocator second-mover. In this encounter an individual with the ordering $(C, C) \succ (D, D)$ establishes the (C, C) outcome, while an individual with $(D, D) \succ (C, C)$ brings about the (D, D) outcome. Thus the ordering with $(C, C) \succ (D, D)$ weakly dominates the other ordering in the μ -game. Note that upon ignoring all these weakly dominated preference types, Assumption 2 is automatically satisfied.

We then obtain the following feasible preference types, where R , M , A and P denote preference types with the best reply behavior of the Reciprocator, Materialist and Altruist type:

- A1: $(C, C) \succ (D, C) \succ (C, D) \succ (D, D)$.
- A2: $(C, C) \succ (C, D) \succ (D, D) \succ (D, C)$.
- A3: $(C, D) \succ (C, C) \succ (D, C) \succ (D, D)$.
- R1: $(C, C) \succ (D, C) \succ (D, D) \succ (C, D)$.
- R2: $(C, C) \succ (D, D) \succ (D, C) \succ (C, D)$.
- R3: $(C, C) \succ (D, D) \succ (C, D) \succ (D, C)$.
- M: $(D, C) \succ (C, C) \succ (D, D) \succ (C, D)$.

- $P1: (D, C) \succ (C, C) \succ (C, D) \succ (D, D)$.
- $P2: (D, C) \succ (C, D) \succ (C, C) \succ (D, D)$.
- $P3: (C, D) \succ (D, C) \succ (C, C) \succ (D, D)$.

Once more, given Assumption 4 above, certain relationships hold between these remaining preference types:

Lemma 2 *In the μ -game, (i). The $R1$ and $R2$ preference types are behaviorally equivalent, as are $P1$ and $P2$, and $A1$ and $A2$. (ii). The $R3$ preference type is weakly dominated by $R1$ and $R2$. (iii). The $A2$ and $A3$ preference types are weakly dominated by $A1$. (iv). The $P3$ preference type is weakly dominated by $P1$ and $P2$.*

The proof of (ii) follows from the fact that an $R3$ player, due to her ranking $(C, D) \succ (D, C)$, performs worse than an $R1$ or $R2$ individual against any Paradoxical type. The $R3$ player, when first-mover in the sequential game, bring about outcome (C, D) , giving her money payoff b , while an $R1$ or $R2$ individual will establish (D, C) , giving them money a . (iii) and (iv) hold for exactly the same reason. Note that upon deleting the weakly dominated preference type $R3$, Assumption 3 is satisfied. Deleting $A2$, $A3$ and $P3$ as well, and merging the behaviorally equivalent preference types, only four preference types survive: R (merger of $R1$ and $R2$), M , A (merger of $A1$ and $A2$) and P (from $P1$ and $P2$). Since these give exactly the four pairs of best replies we have postulated from the beginning, we may regard our definition of a preference type as a pair of pure best replies as a legitimate short-cut.

8 Summary

In this paper we attempt to model, in the simplest possible way, how the features of the 'Game of Life' affect the evolutionarily stable preferences that people develop. We expand the set of strategic situations from containing one to containing two strategic situations. We show that the dynamics of preference evolution, and the stability of the endogenous preference distribution, is qualitatively different from the case where players are engaged in a single strategic situation. There is, under certain conditions, a unique asymptotically stable preference distribution where several preferences co-exist. It seems worthwhile and important for a research program seeking to endogenize preferences that the Game of Life, on which preference evolution is defined, is modelled explicitly and carefully. We believe our results may contribute to providing an evolutionary foundation for the experimentally observed fact that many individuals have social preferences that differ from the materialistic preferences that are normally assumed.

9 Appendix

Proof of Proposition 1: The equations giving the expected money payoffs to the A , R and M preference type are as follows (cf. Table 1): $\pi(A, x) = 1 - x_M + x_M b$,

$\pi(R, x) = x_A + x_R\pi_{RR}$ and $\pi(M, x) = x_A a$. Solving the system $\pi(A, x) = \pi(R, x)$ and $\pi(R, x) = \pi(M, x)$, using $1 = x_A + x_R + x_M$, gives the solutions in (2). It is straightforward to verify that we have $0 < x_i^* < 1$ for all $i = A, R, M$ under the stated conditions on the money payoffs and $\pi_{RR} = \lambda$, $0 < \lambda < 1$.

There are four equilibria for the Replicator dynamic: The three vertices $(1, 0, 0)$, $(0, 1, 0)$ and $(0, 0, 1)$, and x^* . Neither of the two first vertices is a Nash equilibrium, hence each is an unstable equilibrium. The vertex $(0, 0, 1)$ is a Nash equilibrium, but still unstable. To see this, suppose that a small perturbation takes $(0, 0, 1)$ to $x' = (0, \epsilon, 1 - \epsilon)$ where $\epsilon > 0$; at x' , the R -type earns a strictly higher expected monetary payoff than M , hence the proportion of R types increases at the expense of M types, contradicting stability. Thus we need to determine the stability of the interior equilibrium x^* . We will first, instead of the matrix in Table 1, study the equivalent matrix, obtained by subtracting the number 1 (1) [b] from all entries in the first (second) [third] column:

	A	R	M
A	0	0	0
R	0	$\pi_{RR} - 1$	$-b$
M	$a - 1$	-1	$-b$

Or, in abbreviated form,

	A	R	M
A	0	0	0
R	α	β	γ
M	δ	ϵ	θ

Denote this matrix by A , with typical element a_{ij} , where $i, j = A, R, M$, and let $x = (x_A, x_R, x_M)$ be a column vector. Hofbauer (1981) showed that the three-dimensional Replicator dynamic, $\dot{x} = x [Ax - x^T Ax]$, is equivalent to the two-dimensional Lotka-Volterra dynamic:

$$\dot{x} = x [a_{MR} + a_{MM}x + a_{MA}y]$$

$$\dot{y} = y [a_{AR} + a_{AM}x + a_{AA}y],$$

where $x = x_R/x_A$ and $y = x_M/x_A$. If (x, y) is an interior equilibrium of the Lotka-Volterra dynamic, then

$$x^* = (x_A^*, x_R^*, x_M^*) = (1/(1 + x + y), x/(1 + x + y), y/(1 + x + y)) \quad (10)$$

is an equilibrium for the Replicator Dynamic. Conversely, if x^* is an interior equilibrium for the Replicator Dynamic, then

$$(x, y) = (x_R^*/x_A^*, x_M^*/x_A^*) \quad (11)$$

is an equilibrium for the Lotka-Volterra dynamic. Moreover, results about the stability of equilibria for the Lotka-Volterra dynamic carry over to the Replicator

dynamic and conversely, via the two transformations given above. See Hofbauer and Sigmund (1998), Section 7.5. Our strategy is to use the transformation (11) on x^* and study the stability of the corresponding equilibrium (x, y) for the Lotka-Volterra dynamic; the equilibrium x^* will then have the same stability properties under the Replicator Dynamic.

On using the transformation (11) on x^* , we compute the following equilibrium for the Lotka-Volterra dynamic:

$$x = \frac{b(1-a)}{-b\pi_{RR}} = \frac{a-1}{\pi_{RR}} \quad \text{and} \quad y = \frac{(1-\pi_{RR})(a-1)}{-b\pi_{RR}}.$$

We may verify that when these expressions are used in (10), we indeed obtain the solutions in (2) in the main text.

In characterizing the stability of (x, y) , we can use the very useful characterization of equilibria given in Bomze (1983). Bomze shows that if $\beta x + \theta y = 0$, then (x, y) is a center for the Lotka-Volterra system. We compute

$$\beta x + \theta y = \frac{-b(\pi_{RR} - 1)(a - 1) - b(1 - \pi_{RR})(a - 1)}{-b\pi_{RR}} = 0.$$

We may therefore conclude that (x, y) is a center for the Lotka-Volterra system. Using the results in Hofbauer (1981) mentioned above, this allows us to conclude that our equilibrium x^* , given in (2), is a center for the Replicator Dynamic. ■

Proof of Proposition 2: Consider a population, x , with only R and A -players. We then have $\pi(R, x) = \pi(A, x) = \pi(x, x) = 1$. Moreover, we have $\pi(M, x) = x_A a + (1 - x_A)(1/2)$, so $\pi(M, x) < \pi(x, x)$ when $x_A < 1/(2a - 1)$. Then x is a symmetric Nash equilibrium. To show that x is also a Neutrally Stable Strategy (NSS), and hence stable for the Replicator Dynamic¹², we must verify that $\pi(x, x') = \pi(x', x')$ for any $x' \neq x$ using strategy A and R . Since $\pi(x, x') = \pi(x', x') = 1$, x is an NSS (but not an ESS). ■

Proof of Proposition 3: The equations giving the expected money payoffs are $\pi(A, x) = x_A + x_R - x_M b$, $\pi(R, x) = x_A + x_R[\mu\pi_{RR} + 1 - \mu] + x_M[(1 - \mu)(1/2)]$ and $\pi(M, x) = x_A a + x_R[(1/2)(1 - \mu)]$, where $\pi_{RR} = \lambda$ and $0 < \lambda < 1$. Equating these expected money payoffs and solving them, using $y_A + y_R + y_M = 1$, yields the solution in Proposition 3 in the main text. Moreover, it again follows that y^* is interior under the parameter restrictions. None of the corner equilibria, $(1, 0, 0)$, $(0, 1, 0)$ and $(0, 0, 1)$, correspond to Nash equilibria, so they are all unstable. As in the proof of Proposition 1, we use the results in Hofbauer (1981) and Bomze (1983) to characterize the dynamic stability of y^* . On subtracting 1 (1) [b] from the first (second) [third] column from the matrix in Section 4 we obtain the equivalent matrix:

	A	R	M
A	0	0	0
R	0	$\mu\pi_{RR} - \mu$	$(1/2)(1 - \mu) - b$
M	$a - 1$	$-(1/2)(1 + \mu)$	$-b$

¹²We refer the reader to e.g. Weibull (1995).

Using the transformation in (11), the equilibrium for the Lotka-Volterra dynamic is

$$x = \frac{y_R^*}{y_A^*} = \frac{2(a-1)[2b-1+\mu]}{\mu^2 + 2b(2\pi_{RR}-1)\mu + 2b-1}$$

$$y = \frac{y_M^*}{y_A^*} = \frac{4\mu(1-a)(1-\pi_{RR})}{\mu^2 + 2b(2\pi_{RR}-1)\mu + 2b-1}.$$

Bomze (1983) proves that if $\beta x + \theta y < 0$, then (x, y) is asymptotically stable. We compute

$$\beta x + \theta y = \frac{2(\pi_{RR}-1)(a-1)\mu(\mu-1)}{\mu^2 + 2b(2\pi_{RR}-1)\mu + 2b-1}. \quad (12)$$

The numerator is strictly positive for any $\pi_{RR} \in (0, 1)$ and $\mu \in (0, 1)$. The denominator is strictly negative whenever $\mu \in (\underline{\mu}, \bar{\mu})$, and $\underline{\mu}, \bar{\mu} = -2b\pi_{RR} + b \pm \sqrt{A}$, where $A = 4b^2\pi_{RR}^2 - 4b^2\pi_{RR} + b^2 - 2b + 1$. Since $\underline{\mu} < 0$, $\mu \in (\underline{\mu}, \bar{\mu})$ will hold for all $\mu \in (0, 1)$ whenever $\bar{\mu} \geq 1$ or $\sqrt{A} \geq b(2\pi_{RR}-1)$. This is satisfied whenever $\pi_{RR} \leq 1/2$. Suppose then $\pi_{RR} > 1/2$. Then $\bar{\mu} \geq 1$ is the same as $(1-b)^2 \geq 1$, which always holds. Thus we may conclude that the denominator in (12) is negative, and consequently $\beta x + \theta y < 0$. Thus (x, y) is asymptotically stable for the Lotka-Volterra dynamic. This implies that y^* is asymptotically stable for the Replicator Dynamic. ■

Proof of Proposition 4: Let x denote any population where individuals are Reciprocators (R) and/or Altruists (A). The average payoff in any such population is $\pi(x, x) = 1$, since all individuals co-operate with each other, both in the simultaneous and in the sequential games. Any population x' , where a proportion x'_M are M types, earn expected payoff

$$\pi(x', x) = 1 - x'_M + x_M [\mu(1 - x_R)a + (1/2)(1 - \mu) \{x_R(1 + b) + (1 - x_R)a\}].$$

Thus we have $\pi(x', x) < \pi(x, x)$ when

$$\mu(1 - x_R)a + (1/2)(1 - \mu) \{x_R(1 + b) + (1 - x_R)a\} < 1.$$

We observe that for any given μ , this will hold for sufficiently large x_R . Indeed, it is sufficient for this to hold for all $\mu \in (0, 1)$ that (i). $(1 - x_R)a < 1$ or, equivalently, $x_R > (a - 1)/a$ and (ii). that $x_R(1 + b - a) < 2 - a$; the last inequality holds whenever $a < 2$. Suppose then $a > 2$; then $x_R(1 + b - a) < 2 - a$ is equivalent to $x_R > (a - 2)/(a - 1 - b)$. Thus whenever $x_R > \max\{(a - 1)/a, (a - 2)/(a - 1 - b)\} \equiv \underline{x}_R$, we have $\pi(x', x) < \pi(x, x)$ for all populations x' , so x is a Nash equilibrium. To show that x is a Neutrally Stable Strategy (Maynard-Smith (1982)) we must show that $\pi(x, x') \geq \pi(x', x')$ for any x' satisfying $\pi(x', x) = \pi(x, x)$. However, since for any such x' we have $\pi(x, x') = \pi(x', x') = 1$, this holds. ■

Proof of Proposition 5: We first establish a sufficient condition for when M weakly dominates P . As was shown in the main text, this will hold whenever the M type performs at least as well as the P type against the R type, i.e., $1 - \mu \geq 2\mu\pi_{PR} + (1 - \mu)(1 + b)$. Re-arranging this yields $\mu(b - 2\pi_{PR}) > b$ or

$$\mu < \frac{-b}{2\pi_{PR} - b}. \quad (13)$$

Now the money payoff π_{PR} can be arbitrarily close to a ; this depends on the R and P type's exact utility numbers. It is possible that, in the mixed Nash equilibrium, the R -type plays C with almost probability one, and that the P type plays D with almost probability one. Thus a sufficient condition for (13) to hold even in this case is obtained by replacing the money payoff π_{PR} with a . This gives the condition in the main text. Next, given that weak dominance holds, we may use the following result: If a pure strategy, i , is weakly dominated by another (mixed) strategy, call it x , then either strategy i approaches extinction over time, or those pure strategies against which x strictly outperforms i , die out (see Weibull (1995), Proposition 3.2). In our case $i = P$ and we may choose x to be the pure strategy M . That is, strategy P is one of those strategies against which x outperforms P . This implies that P eventually dies out. ■

Proof of Proposition 6: It is not difficult to see that there can be no Altruists in any stable preference distribution. For the Altruists always co-operate and the Materialists always defect, so in any population with Altruists the proportion of Materialists would increase, contradicting stability. This then implies that in any stable preference distribution all players defect, i.e., they are Materialists and/or Reciprocators. Moreover, sufficiently many must be Materialists, since otherwise the Reciprocators would choose to co-operate, contradicting stability (full details are available from the authors upon request). ■

10 References

Bester, H. and Güth, W. (1998): "Is altruism evolutionarily stable?", *Journal of Economic Behavior and Organization*, 34, 193-209.

Bolle, F. (2000): "Is altruism evolutionarily stable? And envy and malevolence?", *Journal of Economic Behavior and Organization*, 42, 131-133.

Bomze, I. (1983): "Lotka-Volterra Equation and Replicator Dynamics: A Two-Dimensional Classification", *Biological Cybernetics*, 48, 201-211.

Charness, G. and Rabin, M. (2001): "Understanding Social Preferences with Simple Tests", forthcoming in *Quarterly Journal of Economics*.

Cooper, R., DeJong, D., Forsythe, R. and Ross, J. (1996): "Cooperation without Reputation: Experimental Evidence from Prisoner's Dilemma Games", *Games and Economic Behavior*, 12, 187-218.

Dawes, R. and Thaler, R. (1988): "Cooperation", *Journal of Economic Perspectives*, 2, 187-197.

Dekel, E., Ely, J.C. and Yilankaya, O. (1998): "Evolution of Preferences", working paper, available at <http://www.kellogg.nwu.edu/research/math/JeffElyJeffEly/index.html/>.

Ely, J.C. and Yilankaya, O. (2001): "Nash Equilibrium and the Evolution of Preferences", *Journal of Economic Theory*, 97, 255-272.

Fehr, E. and Fischbacher, U. (2002): "Why Social Preferences Matter - The Impact of Nonselish Motives on Competition, Cooperation, and Incentives", *Economic Journal*, 112, C1-C33.

Fehr, E. and Gächter, S. (2001): "Fairness and Retaliation: The Economics of Reciprocity", *Journal of Economic Perspectives*, 14, 159-181.

Fehr, E. and Gächter, S. (1998): "Reciprocity and economics: The economic implications of *Homo Reciprocans*", *European Economic Review*, 42, 845-859.

Fehr, E. and Schmidt, K. (1999): "A Theory of Fairness, Competition and Cooperation", *Quarterly Journal of Economics*, 114, 817-868.

Fehr, E. and Schmidt, K. (2001): "Theories of Fairness and Reciprocity - Evidence and Economic Applications", forthcoming in: Dewatripont, M., Hansen, L. and Turnovsky, S. (Eds.): *Advances in Economics and Econometrics - 8th World Congress, Econometric Society Monographs*.

Fershtman, C. and Weiss, Y. (1998): "Why do we care what others think about us?", 133-151 in Ben-Ner, A. and Putterman, L. (eds.): *Economics, Values and Organization*, Cambridge University Press.

Frank, R. (1988): *Passions Within Reasons*, W.W. Norton & Co.

Guttman, J. M. (2000): "On the evolutionary stability of preferences for reciprocity", *European Journal of Political Economy*, 16, 31-50.

Güth, W.: (1995): "An Evolutionary Approach to Explaining Cooperative Behavior by Reciprocal Incentives", *International Journal of Game Theory*, 24, 323-344.

Güth, W. and Kliemt, H. (1998): "The indirect evolutionary approach: Bridging the gap between rationality and adaptation", *Rationality and Society*, 10, 377-399.

Güth, W. and Napel, S. (2002): "Inequality Aversion in a Variety of Games - An Indirect Evolutionary Approach", Working paper 23-2002, Max Planck Institute for Research into Economic Systems.

Güth, W., Schmittberger, R. and Schwarz, B. (1982): "An experimental analysis of ultimatum bargaining", *Journal of Economic Behavior and Organization*, 3, 367-388.

Güth, W. and Yaari, M.: (1992): "An Evolutionary Approach to Explain Reciprocal Behavior in a Simple Strategic Game". In Witt, Ulrich (ed.): *Explaining Process and Change - Approaches to Evolutionary Economics*, Ann Arbor, MI: University of Michigan Press.

Hofbauer, J. (1981): "On the occurrence of limit cycles in the Volterra-Lotka equation", *Nonlinear Analysis, Theory, Methods and Applications*, 5, p. 1003-1007.

Hofbauer, J. and Sigmund, K.: *Evolutionary Games and Population Dynamics*, Cambridge University Press, 1998.

Huck, S. and Oechssler, J. (1999): "The Indirect Evolutionary Approach to Explaining Fair Allocations", *Games and Economic Behavior*, 28, 13-24.

Kandori, M. (1992): "Social Norms and Community Enforcement", *Review of Economic Studies*, 59, 63-80.

Ledyard, John O., "Public Goods: A Survey of Experimental Research." In Kagel, John, and Alvin Roth, eds., *Handbook of Experimental Economics*. Princeton: Princeton University Press, 1995, 111-194.

Levine, D. (1998): "Modeling Altruism and Spitefulness in Experiments", *Review of Economic Dynamics*, 1, 593-622.

Maynard-Smith, J. (1982): *Evolution and the Theory of Games*, Cambridge University Press.

Ok, E. and Vega-Redondo, F. (2001): "On the Evolution of Individualistic Preferences: An Incomplete Information Scenario", *Journal of Economic Theory*, 97, 231-254.

Possajennikov, A. (2000): "On the evolutionary stability of altruistic and spiteful preferences", *Journal of Economic Behavior and Organization*, 42, 125-129.

Possajennikov, A. (2002): "Two-Speed Evolution of Strategies and Preferences in Symmetric Games", SFB 504 discussion paper 02-03, University of Mannheim.

Poulsen, A. (2001): "Reciprocity, Materialism and Welfare: An Evolutionary Model", Working Paper 01-3, Department of Economics, Aarhus School of Business, Denmark.

Poulsen, A. (2003): "Altruism and Welfare when Preferences are Endogenous", unpublished manuscript, Department of Economics, Aarhus School of Business.

Rabin, M. (1993): "Incorporating Fairness into Game Theory and Economics", *American Economic Review*, 83, 1281-1302.

Robson, A. J. (1990): "Efficiency in evolutionary games: Darwin, Nash and the secret handshake", *Journal of Theoretical Biology*, 144, 379-396.

Roth, A. (1995): Bargaining experiments. In Kagel, J. and Roth, A. (eds): *Handbook of Experimental Economics*, 253-348, Princeton University Press.

Samuelson, L. (2001): "Analogies, Adaptation and Anomalies", *Journal of Economic Theory*, 97, 320-366.

Sandholm, W. (2001): "Preference Evolution, Two-Speed Dynamics and Rapid Social Change", *Review of Economic Dynamics*, 4, 637-679.

Sethi, R. and Somanathan, E. (2001): "Preference Evolution and Reciprocity", *Journal of Economic Theory*, 97, 273-297.

Sethi, R. and Somanathan, E. (2003): "Understanding Reciprocity", *Journal of Economic Behavior and Organization*, 50, 1-27.

Sobel, J. (2001): "Interdependent Preferences and Reciprocity", working paper, Department of Economics, University of California, San Diego.

Taylor, P.D. and Jonker, L.B.(1978): "Evolutionarily Stable Strategies and Game Dynamics", *Mathematical Biosciences*, 40, 145-156.

Weibull, J.W. (1995): *Evolutionary Game Theory*, MIT Press.

Zeeman, E.C. (1980): "Dynamics of Evolution of Animal Conflicts", *Journal of Theoretical Biology*, 89, 249-270.