

From reciprocity to unconditional altruism through signalling benefits

Arnon Lotem^{*}, Michael A. Fishman and Lewi Stone

Department of Zoology, Faculty of Life Sciences, Tel-Aviv University, Tel-Aviv 69978, Israel

Cooperation among genetically unrelated individuals is commonly explained by the potential for future reciprocity or by the risk of being punished by group members. However, unconditional altruism is more difficult to explain. We demonstrate that unconditional altruism can evolve as a costly signal of individual quality (i.e. a handicap) as a consequence of reciprocal altruism. This is because the emergent correlation between altruism and individual quality in reciprocity games can facilitate the use of altruism as a quality indicator in a much wider context, outside the reciprocity game, thus affecting its further evolution through signalling benefits. Our model, based on multitype evolutionary game theory shows that, when the additive signalling benefit of donating help exceeds the cost for only some individuals (of high-quality state) but not for others (of low-quality state), the population possesses an evolutionarily stable strategy (ESS) profile wherein high-quality individuals cooperate unconditionally while low-quality individuals defect or play tit-for-tat (TfT). Hence, as predicted by Zahavi's handicap model, signalling benefits of altruistic acts can establish a stable generosity by high-quality individuals that no longer depends on the probability of future reciprocation or punishment.

Keywords: reciprocal altruism; evolutionary game theory; handicap principle; signalling; cooperation

1. INTRODUCTION

The evolution of altruism and cooperation is still one of the most challenging problems in evolutionary theory that continues to attract empiricists and theoreticians alike (Sober & Wilson 1998; Reeve 2000; Wedekind & Milinski 2000; Clutton-Brock *et al.* 2001; Fehr & Gächter 2002; Milinski *et al.* 2002). Altruism among genetically related individuals can evolve by kin selection (Hamilton 1964), but unconditional help offered to unrelated individuals is more difficult to explain. Some degree of cooperation among non-kin may be explained by direct or indirect reciprocity (Trivers 1971; Brems 1996; Nowak & Sigmund 1998a; Lotem *et al.* 1999a; Wedekind & Milinski 2000; Leimar & Hamerstein 2001), or by altruistic punishment (Gintis 2000; Fehr & Gächter 2002). Yet, these models are unlikely to account for unconditional altruism in which help may be given to poor and distant individuals who are unable to respond by either reward or punishment. A different and relatively neglected theory, which may account for unconditional altruism, was proposed by Zahavi (1977, 1987, 1995), based on his handicap principle (Zahavi 1975, 1987). According to Zahavi, seemingly altruistic acts are actually costly signals of quality (i.e. handicaps) by which the performer gains social prestige. For example, by advertising its quality to group members through the performance of costly helping behaviours, a helper might gain direct benefits, such as attracting mates or deterring competitors. Interestingly, however, despite the fact that the handicap principle became a mainstream approach in the study of communication and sexual selection (Grafen 1990; Maynard Smith & Harper 1995; Hasson 1997; Johnston 1997; Lachmann *et al.* 2001), a first attempt to model the idea

of altruism as a handicap has been made only recently (Gintis *et al.* 2001; see § 4).

Zahavi's idea has been mentioned in several recent theoretical studies in which image and reputation were found to be critical elements in forming stable cooperation (indirect reciprocity: Nowak & Sigmund 1998a; Milinski *et al.* 2002; ultimatum games: Nowak *et al.* 2000; and standing strategy: Leimar & Hamerstein 2001). However, in contrast to Zahavi's idea, these models make the simplifying assumption of a uniform population with no individual variation in quality. Thus, the image or reputation being considered can only transmit information on past behaviours but not on individual quality. The first attempt to analyse how individual variation in quality affects reciprocity games may be attributed to Boyd (1992). However, it was Leimar (1997) who was the first to recognize that since the level of help provided by an individual depends on its own quality, it can also transmit information about individual quality, which, in turn, can be used for partner choice within the reciprocity game. A very similar argument, this time with a clear reference to the handicap principle, has been made by Roberts (1998). Roberts' verbal argument took the idea even further, by suggesting that competition for high-quality partners should result in competitive altruism with an escalation in generosity that is rewarded by the value of altruism as a signal of high quality. The next claim made by both Roberts (1998) and Lotem *et al.* (1999b) was to point out that individuals are also likely to use information gathered during reciprocal interaction when they meet the same individuals in other social contexts, such as mate choice or competition over resources. This possibility implies that the consequences of cooperative strategies should be considered in a much wider context than that of a specific reciprocity game. If the image, or the reputation, gained by employing certain strategies during reciprocal interactions can also entail additive signalling benefits due to their effect in other

^{*} Author for correspondence (lotem@post.tau.ac.il).

social contexts, it is almost inevitable that these benefits play some part in shaping these cooperative strategies. According to this intuition, models of reciprocal altruism should probably be extended to consider such signalling benefits in the players' payoff matrix.

The potential effect of signalling benefits on reciprocal altruism is not merely a remote speculation. Several recent game-theory models and computer simulations have analysed the stabilizing effect of individual variation in quality on reciprocal altruism (Leimar 1997; Lotem *et al.* 1999a; Fishman *et al.* 2001; Sherratt & Roberts 2001). An important feature in all of these models is that, at equilibrium, reciprocal help is predominantly being performed by individuals in good phenotypic condition (high-quality state) while poor phenotypes (low-quality state) tend to defect. Hence, altruistic behaviour is correlated with phenotypic quality, and this correlation is already guaranteed to be independent of the potential further role of signalling. Under these conditions, using the level of altruism to assess quality is immediately beneficial provided one needs to acquire such information. Since information on individual quality is known to affect a wide range of social interactions (mainly the outcomes of fights, mate choice and partner choice), the first mutant that uses altruism as a source of information is expected to spread successfully, resulting in a population in which altruism carries additive signalling benefits to the altruist.

We extend the method of *multitype evolutionary game theory* (Cressman 1992) to the case of *frequency-dependent payoff matrices*, and analyse the potential effect of such signalling benefits on reciprocal altruism. We consider signalling benefits in a matrix model previously used to analyse the effect of individual variation in the quality on reciprocal interactions (Fishman *et al.* 2001). This previous analysis showed that phenotypic variations in quality results in an evolutionarily stable strategy (ESS) profile where high-quality individuals play tit-for-tat (TfT) while individuals in poor quality, which cannot afford reciprocity, defect. Here, we show that the introduction of signalling benefits into this model can result in an ESS profile where high-quality individuals (that used to play TfT in the absence of signalling benefits) are cooperating unconditionally, while low-quality individuals defect or play TfT. In other words, we demonstrate that signalling benefits of altruistic acts can establish a stable unconditional altruism by high-quality individuals, which no longer depends on the probability of future reciprocation. Moreover, we show that the conditions for the evolution of this altruism are identical to those of Zahavi's handicap principle: the cost of the signal (of donating favours, in this case) must be of a magnitude that makes it affordable only to some fraction of the population (those in high-quality state) but not to others (those in low-quality state).

2. RECIPROCITY WITHOUT SIGNALLING

We begin by describing the simple case of a symmetric game theoretical model, as used in Fishman *et al.* (2001), and use this as a platform for introducing variation in quality and signalling benefits in a more complex multitype model for heterogeneous populations.

(a) *Symmetric model*

Three evolutionary game strategies (heritable behaviour phenotypes) are considered: unconditional altruists (UAs) that help all other individuals indiscriminately; defectors (DEs) that solicit, but never donate help; and conditional altruists, or TfT players, who retaliate for each defection by refusing help in the future interactions, but otherwise act as UAs. That is, a TfT player always helps UAs and other TfT players, but helps DEs only when it lacks information to classify them. Thus, the response to a request for help by a TfT player depends on its memory of previous interactions.

Denoting the probability that an individual requesting help has been requested to help recently enough for its response to be remembered by $0 < r < 1$, we obtain the following TfT response scheme: with a probability of r , an individual requesting help is correctly classifiable and will be helped if classified as a UA or TfT player, but will be refused help if classified as a DE. With a probability of $1 - r$, an individual requesting help is unclassifiable and will thus be helped (see fig. 1 in Fishman *et al.* 2001). The value of r depends on the probability of repeated interactions and the fidelity of memory and individual recognition.

We assume that individuals meet each other at random and that each individual will eventually interact with many other individuals over its lifetime. We can therefore consider the payoffs for each behavioural strategy in terms of the individual's average accumulated payoffs over a lifetime. Let us denote the average (per capita) accumulated benefits of receiving help over a lifetime by B , and the average lifelong costs of donating help by C (we use capital letters to distinguish these, per lifespan, payoffs from the per encounter payoffs more usual in the literature (cf. Nowak & Sigmund 1998b)).

In these terms, the payoff matrix for both giving and receiving help is

$$P = \begin{array}{ccc} \text{UA} & \text{TfT} & \text{DE} \\ \left(\begin{array}{ccc} B - C & B - C & -C \\ B - C & B - C & -(1 - r)C \\ B & (1 - r)B & 0 \end{array} \right) & \begin{array}{l} \text{UA} \\ \text{TfT} \\ \text{DE} \end{array} \end{array} \quad (2.1)$$

Note that the entry P_{ij} represents the payoffs for an average player of strategy i (horizontal rows) upon interacting with an average player of a strategy j (vertical columns).

As shown by our previous analysis (Fishman *et al.* 2001), independent of whether $C < B$, or whether $C > B$, system (2.1) has a unique ESS solution, DE, which is: DEs displace individuals using alternative strategies, resulting in a population consisting of DEs only.

(b) *A multitype model for heterogeneous populations*

We now present the more realistic case of heterogeneous population that is divided into two classes: *low-quality* individuals for whom costs of reciprocity exceed its benefits, versus *high-quality* individuals for whom reciprocity yields net benefits. The membership in a class is not necessarily hereditary—a reader might find it convenient to think of these quality classes as juveniles and mature

individuals, respectively, or as individuals in poor or good physiological conditions. We shall denote the frequency of low-quality individuals by $0 < q < 1$ (cases $q = 0, 1$ have been addressed in the previous section). We retain the use of B for the accumulated lifelong benefits of reciprocity. Similarly, we retain C as the average accumulated lifelong costs of altruism in the high-quality class, and where $C < B$ since the benefits of reciprocity exceed the costs. For the low-quality class, D is taken as the average accumulated lifelong costs of altruism, where, by definition, $B < D$. Putting these relationships together gives $C < B < D$. Using r as in § 2a, and using the subscripts H and L to denote the (high and low, respectively) quality classes, we have the following payoff matrices: P_{HH} , P_{HL} , P_{LH} , P_{LL} , where the first subscript defines the focal (recipient of the payoff, the row strategy) and the second subscript defines the opponent.

$$P_{HH} = (1 - q) \begin{matrix} \text{UA}_H & \text{TfT}_H & \text{DE}_H \\ \begin{pmatrix} B - C & B - C & -C \\ B - C & B - C & -(1 - r)C \\ B & (1 - r)B & 0 \end{pmatrix} & \begin{matrix} \text{UA}_H \\ \text{TfT}_H \\ \text{DE}_H \end{matrix} \end{matrix}, \quad (2.2)$$

$$P_{HL} = q \begin{matrix} \text{UA}_L & \text{TfT}_L & \text{DE}_L \\ \begin{pmatrix} B - C & B - C & -C \\ B - C & B - C & -(1 - r)C \\ B & (1 - r)B & 0 \end{pmatrix} & \begin{matrix} \text{UA}_H \\ \text{TfT}_H \\ \text{DE}_H \end{matrix} \end{matrix}, \quad (2.3)$$

$$P_{LH} = (1 - q) \begin{matrix} \text{UA}_H & \text{TfT}_H & \text{DE}_H \\ \begin{pmatrix} B - D & B - D & -D \\ B - D & B - D & -(1 - r)D \\ B & (1 - r)B & 0 \end{pmatrix} & \begin{matrix} \text{UA}_L \\ \text{TfT}_L \\ \text{DE}_L \end{matrix} \end{matrix}, \quad (2.4)$$

$$P_{LL} = q \begin{matrix} \text{UA}_L & \text{TfT}_L & \text{DE}_L \\ \begin{pmatrix} B - D & B - D & -D \\ B - D & B - D & -(1 - r)D \\ B & (1 - r)B & 0 \end{pmatrix} & \begin{matrix} \text{UA}_L \\ \text{TfT}_L \\ \text{DE}_L \end{matrix} \end{matrix}. \quad (2.5)$$

Note that the payoffs depend on the fixed frequencies of the two quality types in the population. For example, every element of the P_{HH} and P_{LH} is multiplied by $(1 - q)$ because this is the probability to encounter a high-quality opponent.

The analysis of the above system (Fishman *et al.* 2001) showed that as long as r is sufficiently high (i.e. $r > C/B$) and q , the proportion of individual in low-quality state is not too high (i.e. $q < (rB - C)/r(B - C)$), the population can achieve the cooperative ESS profile (TfT_H , DE_L), that is, play TfT when in a ‘high-quality’ state and defect when in a state of ‘low-quality’.

3. INTRODUCING SIGNALLING BENEFITS

Our reciprocal cooperation results indicate that individuals cooperate when in a high-quality state and defect

when their quality state is low (see also Sherratt & Roberts 2001). Thus, when comparing individuals within a population, the frequency of providing help will be closely correlated with average quality. In social animals, if obtaining information on others’ quality is adaptive for numerous reasons (mainly for decisions associated with mate choice, competitive interactions and partner choice), the first mutant to use helping levels to assess quality in these other contexts will immediately gain a fitness advantage. This is especially feasible since the detectability and memorability required for signal evolution (Guilford & Dawkins 1991; Johnston 1997) are already fulfilled by the requirement for classifying and remembering players’ behaviour in the reciprocity game (i.e. $r > C/B$, see § 2b). In other words, individuals in the reciprocity game already accumulate information on helping behaviour. The only step necessary is to also use the same stored information when making decisions in other contexts, outside the reciprocity game. When this adaptive mutation increases in frequency, and eventually reaches fixation, helping is also beneficial because it signals high quality and can therefore attract mates or deter competitors (Reyer 1986; Putland 2001; Lotem *et al.* 1999b).

Following this logic, we can now introduce signalling benefits to helping behaviour in our reciprocity model. We assume that, each time an individual is providing help, he gains some added signalling benefit as a result of advertising his being at a high-quality state (or more precisely, he increases the probability that this will be noticed and interpreted as a signal of quality—thereby yielding a signalling benefit). It is important to note that over a lifetime, unconditional altruists will accumulate more signalling benefits than conditional altruists (TfT players) because the former always provide help when requested, while the latter may refuse to help DEs. These differences in payoffs between UA and TfT players depend on the frequency of defection behaviours in the population and can only be modelled if the payoff matrix itself is frequency dependent.

We denote the maximum *additive signalling benefits* by S . Let the frequencies of UA_H , TfT_H and DE_H be x_1 , x_2 , x_3 and the frequencies of UA_L , TfT_L and DE_L be y_1 , y_2 , y_3 . And let us denote the ratio of the benefits to the Conditional Altruists (TfT players) to the benefits for UAs by $\psi(\mathbf{x}, \mathbf{y})$, where $0 < \psi(\mathbf{x}, \mathbf{y}) < 1$.

In these terms, the payoff matrices for high-quality players are given by

$$\frac{P_{HH}}{1 - q} = \begin{matrix} \text{UA}_H & \text{TfT}_H & \text{DE}_H \\ \begin{pmatrix} B - C + S & B - C + S & S - C \\ B - C + \psi(\mathbf{x}, \mathbf{y})S & B - C + \psi(\mathbf{x}, \mathbf{y})S & \psi(\mathbf{x}, \mathbf{y})S - (1 - r)C \\ B & (1 - r)B & 0 \end{pmatrix} & \begin{matrix} \text{UA}_H \\ \text{TfT}_H \\ \text{DE}_H \end{matrix} \end{matrix}, \quad (3.1)$$

$$\frac{P_{HL}}{q} = \begin{matrix} \text{UA}_L & \text{TfT}_L & \text{DE}_L \\ \begin{pmatrix} B - C + S & B - C + S & S - C \\ B - C + \psi(\mathbf{x}, \mathbf{y})S & B - C + \psi(\mathbf{x}, \mathbf{y})S & \psi(\mathbf{x}, \mathbf{y})S - (1 - r)C \\ B & (1 - r)B & 0 \end{pmatrix} & \begin{matrix} \text{UA}_H \\ \text{TfT}_H \\ \text{DE}_H \end{matrix} \end{matrix}. \quad (3.2)$$

The payoff matrices for low-quality players are given by

$$\frac{P_{UH}}{1-q} = \begin{pmatrix} UA_H & TtT_H & DE_H \\ B-D+S & B-D+S & S-D \\ B-D+\psi(x,y)S & B-D+\psi(x,y)S & \psi(x,y)S-(1-r)D \\ B & (1-r)B & 0 \end{pmatrix} \begin{matrix} UA_L \\ TtT_L \\ DE_L \end{matrix} \tag{3.3}$$

$$\frac{P_{LL}}{q} = \begin{pmatrix} UA_L & TtT_L & DE_L \\ B-D+S & B-D+S & S-D \\ B-D+\psi(x,y)S & B-D+\psi(x,y)S & \psi(x,y)S-(1-r)D \\ B & (1-r)B & 0 \end{pmatrix} \begin{matrix} UA_L \\ TtT_L \\ DE_L \end{matrix} \tag{3.4}$$

We derive $\psi(x,y)$ as follows. Assume that additive signalling benefits increase linearly with the amount of helping (since helping is also signalling). Now let us consider n interactions in which an individual is approached for help. The probability that a UA helps is unity. Since a TtT player does not help DEs, the probability that a TtT player helps a random applicant is $p = 1 - r[(1 - q)x_3 + qy_3]$. Since the distribution is binomial (either helps or does not), the expectations for n interactions are $E_{UA} = n$ and $E_{TtT} = np$, respectively. Hence

$$\psi(x,y) \equiv E_{TtT}/E_{UA} = p = 1 - r[(1 - q)x_3 + qy_3] \geq 1 - r > 0. \tag{3.5}$$

Unlike the conventional evolutionary games, where the elements of the payoff matrix are constant, system (3.1)–(3.4) is a *multitype evolutionary game with frequency-dependent payoff matrices*. At present, the general theory for analysing such games is still to be formulated. However, in electronic Appendix A (available on The Royal Society’s Publications Web site), we show that the specific case of system (3.1)–(3.4) satisfies the constraints on the evolutionary stability criteria for two-player games with constant payoff matrices that was developed by Cressman (1992), and can therefore be analysed. In electronic Appendix B, we show that system (3.1)–(3.4) has six ESS solutions and one *evolutionary stable set* (ES set) solution.

While the specific conditions for each solution are detailed in electronic Appendix B, the main results can be summarized as follows:

ESS profile (\oplus denotes a mixed solution)	payoff constraints	ESS is attained for:
1. (UA _H , DE _L)	$C < S < D$	all r and q
2. (UA _H , TtT _L \oplus DE _L)	$C < S < D$	$r > \rho, q > \theta_3$
3. (TtT _H , DE _L)	$S < C$	$r > C/B, q < \theta_2$
4. (DE _H , DE _L)	$S < C$	all r and q
5. (TtT _H , TtT _L \oplus DE _L)	$S < C$	$r > \rho, q > \theta_1$
6. (TtT _H \oplus DE _H , DE _L)	$S < C$	all r , and $q < \theta_2$
7. ES set: (UA _H , UA _L)	$S > D$,	

where

$$\rho = \frac{D-S}{B}; \theta_1 = \frac{r-\rho}{r(1-\rho)}; \theta_2 = \frac{rB-C+S}{r(B-C+S)}; \theta_3 = \frac{\rho}{r}.$$

We see that when S is intermediate ($C < S < D$), the system can reach two possible ESS solutions (UA_H, DE_L) and (UA_H, TtT_L \oplus DE_L) (i.e. nos. 1 and 2) in which individuals in a high-quality state help unconditionally. These results make intuitive sense because when the signalling benefit S is greater than the cost of helping (C), S alone should make helping profitable even without future reciprocity. This intuition is confirmed in the first ESS solution (no. 1) for which the stability of the unconditional altruism is attained for all r and q (i.e. independent of the parameter constraints of reciprocity). At this ESS, altruism is a handicap as predicted by Zahavi (1987): it is a costly signal of quality that yields a net benefit of $(S - C)$ to an honest signaller and a fitness loss of $(S - D)$ to a potential cheater (i.e. to an individual of poor quality who tries to signal high quality). It should be noted, however, that although theoretically this ESS is attained for all q , q must be greater than zero to make the condition $C < S < D$ meaningful. If there are absolutely no individuals at a low-quality state for which the cost of providing help is D , helping would be adaptive for the entire population. The system would then behave in practice as in the case of $S > D$, resulting in the ES set solution (UA_H, UA_L) (no. 7) under which every individual helps unconditionally. Realistically, however, at this point helping is no longer correlated with quality and its use as a quality indicator will be selected against, causing S to disappear and the model’s assumption to be invalid. This possible collapse of a signalling system when the signal becomes cheap to everyone was predicted by Zahavi as part of his theory (Zahavi 1987). Note that for all other cooperative solutions (nos. 1–3, 5 and 6) altruism is always positively correlated with quality, in that all ESS profiles contain altruists in high-quality class and DEs in the low-quality class, whereas it is rare to find the converse. This is consistent with the initial assumptions we used to justify the existence of S .

We can see that when $S < C$ unconditional altruism never evolves. However, signalling benefits still play a part, favouring some level of cooperation via a mixed strategy of TtT and defection, even for the low-quality state (solution nos. 2 and 5). It seems that the combined benefit from both signalling and reciprocity can, under some circumstances (i.e. sufficiently high r and q), exceed the cost of reciprocity for low-quality individuals who usually cannot afford it. It is interesting to note that while q , the proportion of low-quality individuals, represents a burden on reciprocity by high-quality individuals (see § 2b; Fishman *et al.* 2001), its effect on low-quality individuals is exactly the opposite. In the present model, low-quality individuals may exhibit some level of TtT only when their frequency in the population (q) exceeds a critical level ($q > (D - S)/Br$ for the case of solution no. 2, and $q > (r - \rho)/r(1 - \rho)$ for solution no. 5).

The model also provides another solution (no. 6), namely (TtT_H \oplus DE_H, DE_L), in which TtT is only part of a mixed strategy when in high quality, but can nevertheless be reached for all values of r (provided that $S > B - C$, $C > B/2$; see electronic Appendix B). This additional scope for TtT is not possible in the absence of signalling, when r must be greater than C/B (Fishman *et al.* 2001), but may be reached if S is sufficiently high to compensate

for the cost of helping many unclassifiable DEs (a result of low r).

Finally, as long as $S < C$, an all-defection ESS (DE_H , DE_L) (no. 4) is also possible depending on initial conditions, and independent of r and q .

4. DISCUSSION

Reciprocity models can explain the evolution of mutual cooperation among individuals but cannot explain the evolutionary stability of unconditional altruism. We analysed the potential role of added signalling benefits in reciprocity games and showed that a stable unconditional altruism can evolve to signal individual quality (i.e. as a handicap).

For concreteness, we based our analysis on extending a previous model in which TfT strategy was taken to represent a typical discriminating strategy in reciprocity models (Fishman *et al.* 2001). Obviously, TfT is not the only possible strategy and reciprocity may be achieved by a wide variety of conditional cooperation strategies (Brembs 1996). Exploring the effect of signalling benefits in reciprocity games involving other conditional strategies is clearly a desirable direction for further research. However, since cooperation is always costly and thus likely to be correlated with individual quality, signalling benefits can potentially affect most other strategies in similar ways, leading to higher levels of cooperation and possibly to unconditional altruism as a handicap.

Our analysis shares some similarities with the work of Gintis *et al.* (2001) on costly signalling and cooperation. These authors modelled the evolution of costly signals in a non-cooperative non-repeated game, and then considered the plausibility that such costly signals will become cooperative. By contrast, our work explored a very different evolutionary path and game structure. We started our analysis from a system that is based on strict reciprocity, and showed that a reciprocity game with repeated interaction can evolve to a point where unconditional 'altruism as a signal' is stable.

A crucial aspect of our analysis is that unconditional altruism is an ESS only when the cost of helping is higher than the benefit of signalling for some individuals (i.e. $D > S > C$). We mentioned in § 3 that this is also a basic condition for the handicap principle, and that Zahavi (1987) has suggested that if an 'inflation process' will make a signal affordable to everyone, the signal will become extinct. In our analysis, this case is illustrated when $S > D$ (ES profile no. 7). At this point, unconditional cooperation should be adapted by everyone. If such a situation does eventuate, the correlation with quality will effectively disappear and the use of altruism as a signal becomes pointless. In theory, this should return the system to its non-signalling state. Realistically, however, this process may not be common. Considering that there are numerous helping opportunities in social animals, it is highly unlikely that all individuals, all of the time, will be in a state where they can afford the cost of providing help. Moreover, when the benefit of signalling is high, individuals of high quality, who compete among themselves, will be selected to invest even more in altruism (Roberts 1998), making these new levels even less affordable to

individuals of a lower quality. In other words, defection by poor phenotypes (Lotem *et al.* 1999a; Fishman *et al.* 2001; Sherratt & Roberts 2001) are unlikely to disappear as a result of signalling benefits because the same signalling benefits may also select for higher levels of help by high-quality individuals. This 'arms-race' perpetuates the affordability gap between quality classes.

Recent work suggests that cooperation among non-kin may be maintained not only by reciprocity, but also by the risk of being punished (Fehr & Gächter 2002). Interestingly, however, in both reciprocity and punishment, cooperation strategies are somehow conditioned (directly or indirectly) by the behavioural strategies played by the recipients of help. The idea of signalling benefits is different in this respect because helping can also be repaid by individuals who merely observe the altruistic acts without ever being potential recipients. In other words, the benefit to the helper depends on the future behaviour of the observers rather than on the future behaviour of the recipients. This may explain, for example, the benefit that a community member may gain by donating money to poor people in a distant undeveloped country, or by funding an agency that provides care to neglected animals. In such cases, the recipients of help are unlikely to reciprocate, or to punish for defection, and help is provided irrespective of the recipients' potential responses. The future response that would matter in this case is that of the donor's community members. Thus, under realistic conditions, 'altruism as a handicap' may still be conditioned, but upon its effect on the audience that assesses the donor's quality rather than by the recipient's response. Therefore, we should not expect that such altruism would really be manifested in the form of a completely unconditional generosity. Instead, we should expect a pattern in which the investment in altruism is like an investment in advertisement. It should depend on its effect on the target audience and on the potential effect of this audience on the fitness of the advertiser.

Signalling benefits may also provide additional insight into recent studies on human cooperation. Milinski *et al.* (2001) found that despite the theoretically predicted superiority of the 'standing strategy' over the mechanism of image scoring (Leimar & Hamerstein 2001), human subjects tend to behave as if they accumulate image rather than satisfying the requirements for good standing. To explain these results Milinski *et al.* (2001) suggested that the standing strategy might be too difficult to apply in terms of memory and cognitive demands. An alternative explanation is that the standing strategy is good and feasible for reciprocity games, but that on a larger scale it is more adaptive to accumulate good image because it can also be used as a signal of high quality under a wide range of additional social contexts. Similarly, signalling benefits may also be involved in explaining the human concept of fairness in ultimatum games (Nowak *et al.* 2000). In these games, contrary to economic rationality, people tend to offer a fair share to other players and to reject low offers that may be perceived as unfair or 'humiliating'. Considering individual variation in quality in such games, the ability to offer a fair share is likely to be correlated with high quality while the tendency to accept low offers may indicate poor individuals who desperately need the money. Thus, fairness may be partly derived from the motivation

to signal that one does not need favours, but can rather afford the cost of a fair game. Finally, even the evidence for altruistic punishment in humans (Fehr & Gächter 2002) may be explained, in part, by signalling benefits. The evolution of punishment in animal societies has usually been associated with dominants who impose their will on subordinates using punishments (Emlen & Wrege 1992; Mulder & Langmore 1993; Clutton-Brock & Parker 1995). In these cases, the punishment acts as an effective negative reinforcement, but the ability to punish can also signal the superior quality of the punisher. Accordingly, the emergence of negative emotions towards DEs as the driving force of altruistic punishment (Fehr & Gächter 2002) may be explained as a general adaptive tendency to retaliate upon losing a resource to another individual in order to signal dominance (i.e. to signal the potential high cost of future conflicts). In this context, it is interesting to recall that in human societies, politicians and policymakers may frequently justify the use of punishments, not by the principle of fairness, but by the strategic argument that a failure to retaliate may be perceived as a signal of weakness.

In conclusion, while the analysis presented may appear complex, the main message of this paper is rather simple. From the moment when we make conventional reciprocity models more realistic by considering individual variation in quality, the interaction with signalling theory becomes almost inevitable and can be theoretically predicted. Altruistic behaviours in nature are therefore expected to be affected by this interaction and their current levels may be maintained, in whole or in part, by signalling benefits.

This work was supported by grants from the Israel Science Foundation (grants 681/96-17.2 and 524/00-17.2) and by the James S. McDonnell Foundation. The authors thank O. Hasson, R. Wagner, R. Johnstone, A. Zahavi, and three anonymous referees for helpful comments and discussions.

REFERENCES

Boyd, R. 1992 The evolution of reciprocity when conditions vary. In *Coalitions and alliances in humans and other animals* (ed. A. H. Harcourt & F. B. M. de Waal), pp. 473–489. Oxford University Press.

Brembs, B. 1996 Chaos, cheating and cooperation: potential solutions to the Prisoner's Dilemma. *Oikos* **76**, 14–24.

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

Clutton-Brock, T. H., Russell, A. F., Sharpe, L. L., Brotherton, P. N. M., McIlrath, G. M., White, S. & Cameron, E. Z. 2001 Effects of helpers on juvenile development and survival in meerkats. *Science* **293**, 2446–2449.

Cressman, R. 1992 *The stability concept of evolutionary game theory*. Berlin: Springer.

Emlen, S. T. & Wrege, P. H. 1992 Parent offspring conflict and the recruitment of helpers among bee-eaters. *Nature* **356**, 331–333.

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

Fishman, M. A., Lotem, A. & Stone, L. 2001 Heterogeneity stabilizes reciprocal altruism interactions. *J. Theor. Biol.* **209**, 87–95.

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

Gintis, H., Smith, E. A. & Bowles, S. 2001 Costly signaling and cooperation. *J. Theor. Biol.* **213**, 103–119.

Grafen, A. 1990 Biological signals as handicaps. *J. Theor. Biol.* **144**, 517–546.

Guilford, T. & Dawkins, M. S. 1991 Receiver psychology and the evolution of animal signals. *Anim. Behav.* **42**, 1–14.

Hamilton, W. D. 1964 The genetic evolution of social behavior. I & II. *J. Theor. Biol.* **7**, 1–52.

Hasson, O. 1997 Towards a general theory of biological signaling. *J. Theor. Biol.* **185**, 139–156.

Johnston, R. A. 1997 The evolution of animal signals. In *Behavioral ecology: an evolutionary approach* (ed. J. R. Krebs & N. B. Davies), pp. 155–178. Oxford: Blackwell Scientific.

Lachmann, M., Számádó, S. & Bergstrom, C. T. 2001 Cost and conflict in animal signals and human language. *Proc. Natl Acad. Sci. USA* **98**, 13 189–13 194.

Leimar, O. 1997 Reciprocity and communication of partner quality. *Proc. R. Soc. Lond. B* **264**, 1209–1215. (DOI 10.1098/rspb.1997.0167.)

Leimar, O. & Hamerstein, P. 2001 Evolution of cooperation through indirect reciprocity. *Proc. R. Soc. Lond. B* **268**, 745–753. (DOI 10.1098/rspb.2000.1573.)

Lotem, A., Fishman, M. A. & Stone, L. 1999a Evolution of cooperation between individuals. *Nature* **400**, 226–227.

Lotem, A., Wagner, R. H. & Balshine-Earn, S. 1999b The overlooked signaling component of nonsignaling behavior. *Behav. Ecol.* **10**, 209–212.

Maynard Smith, J. & Harper, D. G. C. 1995 Animal signals: models and terminology. *J. Theor. Biol.* **177**, 305–311.

Milinski, M., Semmann, D., Bakker, T. C. M. & Krambeck, H.-J. 2001 Cooperation through indirect reciprocity: image scoring or standing strategy? *Proc. R. Soc. Lond. B* **268**, 2495–2501. (DOI 10.1098/rspb.2001.1809.)

Milinski, M., Semmann, D. & Jrambeck, H.-J. 2002 Reputation helps solve the 'tragedy of the commons'. *Nature* **415**, 424–426.

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

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

Nowak, M. A. & Sigmund, K. 1998b The dynamics of indirect reciprocity. *J. Theor. Biol.* **194**, 561–574.

Nowak, M. A., Page, K. M. & Sigmund, K. 2000 Fairness versus reason in the ultimatum game. *Science* **289**, 1773–1775.

Putland, D. 2001 Has sexual selection been overlooked in the study of avian helping behaviour? *Anim. Behav.* **62**, 811–814.

Reeve, H. K. 2000 Unto others—the evolution and psychology of unselfish behaviour by Sober E., Wilson D. S. *Evol. Hum. Behav.* **21**, 65–72.

Reyer, H. U. 1986 Breeder-helper interactions in the pied kingfisher reflects the costs and benefits of cooperative breeding. *Behaviour* **96**, 277–303.

Roberts, G. 1998 Competitive altruism: from reciprocity to the handicap principle. *Proc. R. Soc. Lond. B* **265**, 427–431. (DOI 10.1098/rspb.1998.0312.)

Sherratt, T. N. & Roberts, G. 2001 The role of phenotypic defectors in stabilizing reciprocal altruism. *Behav. Ecol.* **12**, 313–317.

Sober, E. & Wilson, D. S. 1998 *Unto others—the evolution and psychology of unselfish behaviour*. London: Harvard University Press.

Trivers, R. 1971 The evolution of reciprocal altruism. *Q. Rev. Biol.* **46**, 35–56.

Wedekind, C. & Milinski, M. 2000 Cooperation through image scoring in human. *Science* **288**, 850–852.

Zahavi, A. 1975 Mate selection—a selection for a handicap. *J. Theor. Biol.* **53**, 205–214.

Zahavi, A. 1977 Reliability in communication systems and the evolution of altruism. In *Evolutionary ecology* (ed. B. Stonehouse & C. M. Perrins), pp. 253–259. London: Macmillan.

Zahavi, A. 1987 The theory of signal selection and some of its implications. In *Int. Symp. Biological Evolution* (ed. V. P. Delfino), pp. 305–327. Bari: Adriatica Editrice.

Zahavi, A. 1995 Altruism as a handicap: the limitation of kin selection and reciprocity. *J. Avian Biol.* **26**, 1–3.

As this paper exceeds the maximum length normally permitted, the authors have agreed to contribute to production costs.

Visit <http://www.pubs.royalsoc.ac.uk> to see an electronic appendix to this paper.