COMMENTARY

# Reciprocity: you have to think different

R. BOYD

*Department of Anthropology, University of California, Los Angeles, CA, USA and
The Leverhulme Centre for Human Evolutionary Studies, University of
Cambridge, Cambridge, UK*

There has been a proliferation of models seeking to explain the evolution of altruism over the past few years, and Lehmann and Keller should be commended for this very well-reasoned and clearly written attempt to bring order out of chaos. There is much good sense in this paper, especially the emphasis on the importance of assortment for the evolution of altruism, the excellent discussion of green beard models and the important role of constraints, particularly in punishment models. I fear, however, that the elegance and simplicity of their results may mislead some readers about the implications of this work for the evolution of reciprocity.

Lehmann and Keller consider the evolution of reciprocity in between pairs of individuals. Unlike most authors, they assume that helping is a continuous variable. The incremental cost of helping is $C$ and the benefit $B$, so a focal individual who helps an amount $z$ suffers a cost $Cz$ and his partner receives a benefit $Bz$. The strategy set is limited to linear reactive strategies. If a focal's partner helped an amount $z$ on the previous turn, then the focal helps $m\beta z$ on this turn. Then, Lehman and Keller show that altruism among unrelated individuals can be favoured by selection when:

$$\omega m\beta B > C \qquad (1)$$

where $\omega$ is the probability of continued interaction. Notice that if $m\beta = 1$, this model becomes the continuous version of Tit-for-tat (do whatever your partner did on the last turn) and, satisfyingly, the condition for increase is exactly that derived by Axelrod & Hamilton (1981) in the discrete case. This expression seems to suggest that kin selection and reciprocity are similar – in both cases, altruism requires assortative interaction. In kin selection, it is cues of relatedness that allow individuals to assort, and in reciprocity, it is past behaviour. So far so good. However, it is then easy to conclude that if Eqn 1 is satisfied, that reciprocity is a likely evolutionary outcome. If the parallel condition for altruism among relatives is satisfied, genes that give rise to altruism can increase when rare and will also spread to fixation under the usual assumptions of additive gene effects and weak

*Correspondence:* R. Boyd, Department of Anthropology, UCLA, Los Angeles, CA, USA and The Leverhulme Centre for Human Evolutionary Studies, University of Cambridge, Cambridge, UK.
Tel.: 01 310 206 8008; fax: 01 310 206 7833; e-mail: robert.t.boyd@gmail.com

selection – altruism *is* the likely evolutionary outcome. However, this is definitely *not* the case for reciprocity – when Eqn 1 is satisfied, reciprocity can be an evolutionarily stable strategy (ESS), but so do, can every other pattern of behaviour.

When there are repeated social interactions, and interacting individuals can have substantial effects on each other's fitness, virtually any pattern of behaviour can be a stable (or almost stable) evolutionary equilibrium. Game theorists call this result the 'folk theorem' because it was widely known among game theorists in the 1950s, but nobody was exactly sure who first proved it, and complete proofs were not published until much later (e.g. Rubinstein, 1979; Fudenberg & Maskin, 1986). The basic logic of the folk theorem is simple. Suppose there is a strategy that takes the form: do $x$ where $x$ is some behaviour, say alternating cooperate and defect, as long as the other guy does $x$. If the other guy does something else, defect forever. Once a strategy like this becomes common in a population, the only smart thing to do is to do $x$, otherwise you will be punished by defection for as long as the interaction lasts. If interactions go on long enough, the costs of such punishment will exceed the short-run benefits of doing something other than $x$. In short, repeated interactions create the possibility of sanctions, and any behaviour that enough sanctioners are willing to sanction is an equilibrium. For the most part, the logic of the folk theorem applies to evolutionary theory, although there is a subtle and important difference that affects the stability of punishment. The bottom line is that when everything is an equilibrium, showing that reciprocity is an equilibrium too does not really tell you much.

It will be a lot easier to see how big a problem this is if we add just a touch more realism to Lehmann and Keller's model. Consider a situation in which the costs and benefits of cooperation vary from period to period. For example, a first type of interaction might be coalitional aid, a second could be grooming, and a third sharing food. To model this, suppose that there are $n$ types of interactions, and that type $i$ has benefit $b_i$, cost $c_i$, and occurs with probability $p_i$. Strategies specify that an individual must cooperate when his partner is in good standing or he is not. Individuals maintain good standing only as long as they cooperate in a specified subset of interaction types.[1] For example, if there were three interaction types, there would be eight possible strategies: {Ø} {1} {2} {3} {1,2} {1,3} {2,3} {1,2,3}. In each case, individuals cooperate and expect cooperation only

---

[1]Lehmann and Keller restrict the attention to linear reactive strategies. This constrains past behaviour to be a predictor of future behaviour and thus fits nicely into their conceptual scheme. However, there is no reason to think that such strategies are evolutionarily robust. In discrete models, highly nonlinear standing-based strategies like Contrite Tit-for-tat typically out-compete linear reactive strategies like Tit-for-tat or Tit-for-two-tats (Boyd, 1989; Wu & Axelrod, 1995; Leimar, 1997).

for interaction types in the set. Obviously, there are a very large number of strategies if there are very many interaction types. Now, here is the interesting bit. If interactions go on long enough, *any* strategy that produces a net benefit can be an ESS, even strategies that specify cooperation in interaction types in which costs exceed the benefits, i.e. $b_i < c_i$ (Boyd, 1992)! You might think that mutants that do not cooperate in such unproductive interactions would do better, but they do not. The reason is that once a particular strategy becomes common, you have to follow that strategy. If you do not, others will stop cooperating and you will lose the long-run benefits of cooperation. Thus, the theory of reciprocal altruism predicts a vast range of alternative outcomes are possible, which is to say, absent some account of which equilibria are likely, it predicts almost nothing.

For repeated interactions between pairs of individuals, there is a partial solution to this dilemma. When pairs interact, small amounts of relatedness destabilize relatively less cooperative equilibria, but not relatively more cooperative ones (Axelrod & Hamilton, 1981). When individuals are related, then individuals with rare, invading strategies have some chance of interacting. Thus, if the invading strategy produces a big enough mutual benefit, rare interactions between invading pairs can compensate for the low pay-off that invaders achieve in their much more common interactions with dominant type and, as a result, invaders have higher fitness averaged over all interactions than do dominants. Thus, as Axelrod and Hamilton put it, social evolution has a ratchet. The action of the ratchet depends on the relative payoff of the strategies involved. Only small amounts of relatedness may be necessary to allow reciprocating strategies to invade unconditional defection because two defectors often have much lower fitness than two cooperators. Larger amounts of relatedness may be necessary to allow a better reciprocating strategy to invade a poorer one if the differences in average fitness are smaller.

However, the ratchet does not work at all when the benefits of cooperation flow to sizable groups (Boyd & Richerson, 1988). Lehmann and Keller restrict their analysis to interactions between pairs of individuals. This assumption makes little difference for kin selection or green beard models. However, it makes a huge difference in the case of reciprocity (Axelrod & Dion, 1988; Boyd & Richerson, 1988). To see why, suppose that individuals live in groups, and each helping act benefits all group members. For example, the helping behaviour could be an alarm cry that warns group members of an approaching predator, but makes the callers conspicuous and thereby increases their risk of being eaten. Now, consider the fate of a rare defector. If reciprocators use the rule, only cooperate if all others cooperate, presence of this defector causes cooperation in its group to collapse. As long as long-run cooperation pays, cooperators in groups without defectors will have higher fitness, and defectors will not invade. However, if reciprocators use any other more tolerant rule, defectors will get the benefits of cooperation without paying the cost and will invade. Thus, only intolerant reciprocators can persist when common. However, even substantial amounts of relatedness are not enough to allow such intolerant reciprocating strategies to invade when rare. To see why, suppose that interacting individuals are full sibs, and that groups are composed of 10 individuals. Then only one in 512 reciprocators will get any benefit from their attempts to cooperate. The other 511 will suffer the cost of their initial cooperation without any long-run benefit. Lower levels of relatedness or larger groups make things even worse. Thus, unless the long-run benefits are much greater than the short-run costs, reciprocity cannot increase.

You might think that ratchet would work if punishment took some other form – noncooperators can be punished by reduced status, fewer friends and fewer mating opportunities – what Triver's (1971) called 'moralistic reciprocity'. With ordinary reciprocity, the severity of the sanction is limited by the effect of a single individual's cooperation on each other group member, an effect that becomes small as group size increases. Moralistic sanctions can be much more costly to defectors, making it possible for cooperators to induce others to cooperate in large groups even when they are rare. Cowards, deserters and cheaters may be attacked by their erstwhile compatriots and shunned by their society, made the targets of gossip or denied access to territories or mates. Thus, moralistic punishment provides a more plausible mechanism for the maintenance of large-scale cooperation than reciprocity. However, it does not solve the problem of multiple equilibria. In fact, adding relatedness to models with moralistic punishment actually decreases the relative fitness of cooperating strategies (Boyd & Richerson, 1992; Gardner & West, 2004).

Unlike the theory of kin selection, the theory of reciprocal altruism is fundamentally incomplete. If evolutionary change is driven only by individual costs and benefits, then reciprocity and moralistic punishment can stabilize cooperation, but they can also stabilize anything else. The reason is that reciprocity is really a form of coordination. Reciprocating strategies require conformance to a rule, and punish those who deviate. As long as being punished is sufficiently costly, anything can be stabilized. As cooperative behaviours are a tiny subset of all possible behaviours, reciprocity and punishment cannot by themselves explain observed cooperation. Selection will pick out mutually beneficial behaviours only if kin selection or group selection act to favour those strategies that coordinate so as to create mutual benefit.

## References

Axelrod, R. & Dion, D. 1988. The further evolution of cooperation. *Science* **242**: 1385–1390.
Axelrod, R. & Hamilton, W.D. 1981. The evolution of cooperation. *Science* **211**: 1390–1396.

Boyd, R. 1989. Mistakes allow evolutionary stability in the repeated prisoner's dilemma game. *J. Theor. Biol.* **136**: 47–56.

Boyd, R. 1992. The evolution of reciprocity when conditions vary. In: *Coalitions in Humans and Other Animals*, (F. Dewaal & A. H. Harcourt, eds), pp. 473–492. Oxford University Press, Oxford.

Boyd, R. & Richerson, P.J. 1988. The evolution reciprocity in sizable groups. *J. Theor. Biol.* **132**: 337–356.

Boyd, R. & Richerson, P.J. 1992. Punishment allows the evolution of cooperation (or anything else) in sizable groups. *Ethol. Sociobiol.* **13**: 171–195.

Fudenberg, D. & Maskin, E. 1986. The folk theorem for repeated games with discounting and incomplete information. *Econometrica* **54**: 533–554.

Gardner, A. & West, S. 2004. Cooperation and punishment, mainly in humans. *Am. Nat.* **164**: 753–764.

Leimar, O. 1997. Repeated games: a state space approach. *J. Theor. Biol.* **184**: 471–498.

Rubinstein, A. 1979. Equilibrium in supergames with the overtaking criterion. *J. Econ. Theory* **21**: 1–9.

Trivers, R.L. 1971. The evolution of reciprocal altruism. *Q. Rev. Biol.* **46**: 35–57.

Wu, J. & Axelrod, R. 1995. How to cope with noise in the iterated Prisoner's Dilemma. *J. Conflict Resolut.* **39**: 183–189.