

**WORKING PAPER SERIES**

Sergio Beraldo, Robert Sugden

**THE EMERGENCE OF RECIPROCALLY BENEFICIAL COOPERATION**

*Working Paper No. 18/2010*

# The emergence of reciprocally beneficial cooperation

Sergio Beraldo<sup>♦</sup> and Robert Sugden<sup>♠</sup>

July 2010

**Abstract:** This paper offers a new and robust model of the emergence and persistence of cooperation. In the model, interactions are anonymous, the population is well-mixed, and the evolutionary process selects strategies according to material payoffs. The cooperation problem is modelled as a game similar to Prisoner's Dilemma, but there is an outside option of non-participation and the payoff to mutual cooperation is stochastic; with positive probability, this payoff exceeds that from cheating against a cooperator. Under mild conditions, mutually beneficial cooperation occurs in equilibrium. This is possible because the non-participation option holds down the equilibrium frequency of cheating.

**Keywords:** Cooperation; voluntary participation; random payoffs.

**JEL Classification:** C73.

---

<sup>♦</sup> *University of Naples Federico II & ICER.. Dipartimento di Scienze dello Stato, via Mezzocannone 4, 80134, Naples, Italy. s.beraldo@unina.it . S. B. acknowledges financial support from the University of Naples "Federico II" (Programma di breve mobilità) and from the Banco di Napoli Foundation.*

<sup>♠</sup> *University of East Anglia, School of Economics and Centre for Behavioural and Experimental Social Science, Norwich, NR4 7TJ, UK. r.sugden@uea.ac.uk.*

## 1. Introduction

Studies of animal behaviour have found many practices which create collective benefits at some apparent cost or risk to individual participants. Examples include alarm calls, food-sharing, grooming, and participation in inter-group warfare. One of the most fundamental problems in evolutionary biology since Darwin (1859) has been to explain how such forms of cooperation evolve by natural selection. An analogous problem in economics has been to explain how cooperative human practices, such as the fulfilment of market obligations, the provision of public goods and the management of common property resources, are consistent with the traditional assumption of individual self-interest. Many different theories have been proposed by biologists and economists as possible solutions. Among the mechanisms that have been modelled are direct and indirect reciprocity, reputation, third-party punishment, kin selection, group selection, and the ‘green beard’ mechanism.<sup>1</sup>

However, a recent trend in biology has been to question whether such sophisticated explanations are always necessary. Many forms of apparently cooperative behaviour have been found to be forms of mutualism: the ‘cooperating’ individual derives sufficient direct fitness benefit to make the behaviour worthwhile, and any effect on the fitness of others is incidental (e.g. Clutton-Brock, 2002, 2009; Sachs et al., 2004). The Snowdrift game (Sugden, 1986), in which equilibrium involves cooperation by one player and free-riding by the other, is increasingly used in biology as a model of such behaviour. In this paper, we present a new model of the evolution of cooperation which fits with this trend of thought.

The idea that the biological and economic problems of cooperation are isomorphic can be developed in at least two different ways. One approach is to hypothesise that human cooperation in the modern world is a product of genetically hard-wired traits which evolved to equip *homo sapiens* for life in hunter-gatherer societies (e.g. Binmore 1994, 1998). An alternative approach, and the one we will take, is to hypothesise that the emergence and reproduction of human cooperative practices are governed by evolutionary mechanisms that are isomorphic to, but distinct from, those of natural selection. Candidate mechanisms include

---

<sup>1</sup> For an overview of these mechanisms, see Nowak (2006).

trial-and-error learning by individuals, imitation of successful neighbours, and cultural selection through inter-group competition. Analyses which use this approach may be both informed by and informative to theoretical biology. For example, Sugden's (1986) analysis of the emergence of social norms was inspired by the earlier work of theoretical biologists, but it developed new models (in particular, the Snowdrift and Mutual Aid games) which have since been widely used in biology (e.g. Leimar and Hammerstein, 2001; Nowak and Sigmund, 2005). The model that we present in this paper can be interpreted as a representation either of natural selection or of trial-and-error human learning.

Our modelling strategy is distinctive in that it uses three assumptions which in combination rule out most of the mechanisms that feature in existing theories of cooperation. Specifically, we assume that interactions are anonymous, that evolution takes place in a large, well-mixed population, and that the evolutionary process selects strategies according to their material payoffs. The assumption of anonymity excludes mechanisms based on reputation, reciprocity or third-party punishment. The assumption of well-mixedness excludes mechanisms of group or kin selection. The assumption that selection is for material payoffs excludes mechanisms which postulate non-selfish preferences as an explanatory primitive. Working within the constraints imposed by these assumptions, we are able to generate a simple and robust model of cooperation.

Our model adapts the familiar framework of a Prisoner's Dilemma that is played recurrently in a large population. We introduce two additional features, which we suggest can be found in many real-world cases of potentially cooperative interaction, both for humans and for other animals.

The first additional feature is that participation in the game is voluntary. One of the restrictive properties of the Prisoner's Dilemma is that, in any given interaction, an individual must act either *pro*-socially (the strategy of cooperation) or *anti*-socially (the strategy of defection or cheating, which allows a cheater to benefit at the expense of a cooperator). There is no opportunity to be simply *asocial*. We add an *asocial* strategy, that of opting out of the interaction altogether.

The second additional feature is that the payoff that each player receives if they both cooperate is subject to random variation. Before choosing his (or her, or its) strategy, each player knows his own cooperative payoff, but not the other player's. With non-zero probability, the payoff from mutual cooperation is greater than that from cheating against a cooperator. Thus, there are circumstances in which it would be profitable for a player to cooperate if he were sufficiently confident that the other player would cooperate too.

As an illustration of the kind of interaction that our model represents, we offer the following variant of Rousseau's (1755/ 1988, p. 36) story of hunting in a state of nature. Two individuals jointly have the opportunity to invest time and energy to hunt a deer. The hunters can succeed only by acting on a concerted plan out of sight of one another. A hunt begins only if both individuals agree to take part. Each can then cheat by unilaterally pursuing a smaller prey, which the other's deer-hunting tends to flush out and make easier to catch. The anticipated benefit of deer-hunting to an individual, conditional on the other's not cheating, can be different for different individuals and on different occasions. Sometimes, but not always, this benefit is sufficiently low that unilateral cheating pays off.

As a more modern illustration, consider two individuals who make contact through the internet. One of them is offering to sell some good which has to be customised to meet the specific requirements of the buyer; the other is looking to buy such a good. If they agree to trade, each individual invests resources in the transaction (exchanging information, producing and dispatching the good, sending payment). Each may have opportunities to gain by deviating from the terms of the agreement. Sometimes, but not always, the benefit of completing the transaction is sufficiently low that unilateral cheating pays off.

We will show how the interaction of voluntary participation and stochastic payoffs can induce cooperation. Of course, it is well known that voluntary participation can facilitate cooperation when players can distinguish between more and less cooperative opponents. If such distinctions are possible, voluntary participation can allow cooperators to avoid interacting with cheats. This can sustain cooperation without the need for informationally and cognitively more demanding strategies of reciprocity or punishment – an idea that can be traced back to Adam Smith's (1763/ 1978, pp. 538–539) analysis of trustworthiness among

traders in commercial societies. But such mechanisms are ruled out by our anonymity assumption.

In our model, voluntary participation facilitates cooperation by a different route. Because would-be cheats have the alternative option of non-participation, and because non-participation is the best response to cheating, the equilibrium frequency of cheating is subject to an upper limit. If cheating occurs at all, the expected payoff from cheating cannot be less than that from non-participation. Thus, for any given frequency of cooperation, the frequency of cheating is self-limiting. The underlying mechanism is similar to that of the Lotka–Volterra model of interaction between predators and prey: the size of the predator population (the frequency of cheating) is limited by the size of the prey population (the frequency of cooperation).

Clearly, however, this mechanism can support cooperation only if, when the frequency of cheating is sufficiently low, some players choose to cooperate. This could not be the case if, as in the Prisoner’s Dilemma, cooperation was *always* a weakly dominated strategy. In our model, random variation in the payoff from mutual cooperation ensures that players sometimes find it worthwhile to cooperate, despite the risk of meeting a cheat. The players who cooperate are those for whom the benefit of mutual cooperation is sufficient to compensate for this risk. Because cooperators are self-selecting in this way, the average payoff in the game is greater than the payoff to non-participation. In other words, despite the presence of cheats, beneficial cooperation occurs.

In Section 2 we present the model and identify its Nash equilibria. We show that, provided the upper bound of the distribution of cooperative benefit is not too low, the game has at least one equilibrium in which beneficial cooperation occurs. In Section 3 we investigate some comparative-static properties of the model. We show that as the distribution of cooperative benefit becomes more favourable, the maximum frequency of cooperation that is sustainable in equilibrium increases. In Section 4 we examine the dynamics of the model. We show that, in the neighbourhood of equilibria in which some but not all players choose to participate, the dynamics induce cycles similar to those of predator–prey models. In Section 5, we discuss the contribution that our model can make to the explanation of cooperative

behaviour. We show that, despite sharing some features of existing biological models of mutualism and voluntary participation, it isolates a distinct causal mechanism.

## 2. The model: equilibrium properties

We consider a setting with a large number of individuals, interacting anonymously in an indefinitely long sequence of periods. In each period, individuals are randomly matched to play a two-player game. In a representative game between players  $i$  and  $j$ , the benefits from cooperation  $x_i$  and  $x_j$  are independent realizations of a random variable  $X$  whose distribution  $f(\cdot)$  is continuous with support  $[x_{\min}, x_{\max}]$ . Each player knows his own benefit but not that of the other player. Given this knowledge, he chooses one of three options – to cooperate (C), to cheat (D), or not to participate (N). The payoff matrix is shown in Table 1.

Table 1: Payoff matrix for the game

		player $j$		
		N	C	D
player $i$	N	0, 0	0, 0	0, 0
	C	0, 0	$x_i, x_j$	$-b, a$
	D	0, 0	$a, -b$	$-c, -c$

$$x_{\max} > a > x_{\min} \geq 0; b > a > c > 0.$$

The essential features of the game are contained in the structure of best responses. The condition  $x_{\max} > a > x_{\min}$  imposes the assumption that either C or D may be the better response to C, depending on the relevant player's realization of  $X$ . The condition  $b > c$  implies that, as in the Prisoner's Dilemma, D is better than C as a response to D. Given that the payoff to N is normalized to zero,  $a > 0$  implies that cheating gives a higher payoff than non-participation if

the opponent cooperates;  $c > 0$  implies that the opposite is the case if the opponent cheats;  $b > a$  implies that the benefit from cheating a cooperating co-player is less than the cost inflicted on the latter. The condition  $x_{\min} \geq 0$  (which is not essential for our main results) implies that players are never worse off from mutual cooperation than from non-participation.<sup>2</sup>

We assume that, in any given period, for any given player  $i$ , there is some critical value of  $X$  such that  $i$  plays C if and only if  $x_i$  is greater than or equal to this value. Given this assumption, which attributes a minimal degree of rational consistency to players' behaviour, we can represent a strategy for playing the game by two variables – the critical value of  $X$ , denoted by  $\beta$ , and the probability  $\pi$  that D is played, conditional on  $X$  being below that critical value. To simplify the exposition, we impose the harmless condition that  $x_{\min} \leq \beta \leq x_{\max}$ . A strategy  $(\beta, \pi)$  is an *equilibrium* if it is a best reply to itself.

Some significant properties of equilibrium hold for all parameter values. First,  $(\beta, \pi) = (x_{\max}, 0)$  is an equilibrium. In this *non-participation equilibrium*, N is always chosen; players' payoffs are zero, and unilateral deviations lead to neither gain nor loss. Second, there is no equilibrium in which C is played with nonzero probability but D is not played: against an opponent who might play C but will not play D, the best reply chooses D when  $x_i < a$ . Third, there is no equilibrium in which D is played but not C: against an opponent who might play D but will not play C, N is the unique best reply. Thus, only two types of equilibrium participation are possible. Depending on the parameter values, there may be an *interior equilibrium* with  $x_{\min} < \beta < x_{\max}$  and  $0 < \pi < 1$ , characterized by N, D and C being played with nonzero probability; and there may be a *boundary equilibrium* with  $x_{\min} < \beta < x_{\max}$  and  $\pi = 1$ : in this case, D and C are played but not N.

We now analyse these equilibria. Consider any player  $i$  facing an opponent whose strategy is  $(\beta, \pi)$ , in an interaction in which  $x_i = \beta$ . Let  $V_N$ ,  $V_D$ ,  $V_C$  and  $V_M$  be the expected payoffs to player  $i$  from playing N, D, C and M respectively, where M is the mix of D with

---

<sup>2</sup> Provided that this best response structure is maintained, the main implications of the model are preserved. It is not essential that the payoff from playing C against C is stochastic and that all other payoffs are not; but there must be some random variation in the payoffs, such that the best reply to C is sometimes C and sometimes D.



probability  $\pi$  and N with probability  $(1-\pi)$ . Let  $g(x) \equiv F(x)/[1-F(x)]$ , where  $F(\cdot)$  is the cumulative of  $f(\cdot)$ . It is straightforward to derive the following expressions:

$$V_N = 0 \tag{1}$$

$$V_D = [1-F(\beta)]a - F(\beta)\pi c \tag{2}$$

$$V_C = [1-F(\beta)]\beta - F(\beta)\pi b \tag{3}$$

$$V_M = \pi V_D. \tag{4}$$

In analysing equilibrium, it is convenient to work in a  $(\beta, \pi)$  space defined by  $x_{\min} \leq \beta \leq x_{\max}$  and  $\pi \geq 0$ . Notice that this space includes points at which  $\pi > 1$ . Although such points have no interpretation within our model, equations (1)–(4) above define  $V_N$ ,  $V_D$ ,  $V_C$ , and  $V_M$  for *all* values of  $\pi$ . This allows us to define the loci of points in this  $(\beta, \pi)$  space at which the mathematical equations  $V_N = V_D$  and  $V_C = V_M$  are satisfied, and then to characterise equilibria in terms of these loci, imposing the inequality  $\pi \leq 1$  as an additional constraint. This method of analysis is useful in simplifying the proofs of our results.

First, consider the locus of points in the  $(\beta, \pi)$  space at which  $V_N = V_D$ . Any interior equilibrium must be a point on this *ND locus*, with  $0 < \pi < 1$ ; any boundary equilibrium must be a point at which  $V_N \leq V_D$  and  $\pi = 1$ . By (1) and (2), this locus is determined by:

$$V_D = (\text{or } <) V_N \Leftrightarrow a/\pi c = (\text{or } <) g(\beta). \tag{5}$$

This is a continuous and downward-sloping curve which includes the point  $(x_{\max}, 0)$  and is asymptotic to  $\beta = x_{\min}$ . It divides the  $(\beta, \pi)$  space into three regions: the set of points *on* the locus, at which  $V_N = V_D$ ; the set of points *inside* the locus (that is, below and to the left), at which  $V_N < V_D$ ; and the set of points *outside* the locus (that is, above and to the right), at which  $V_N > V_D$ .

Now consider the locus of points at which  $V_C = V_M$ . Every equilibrium must be a point on this *CM locus*, with either  $\pi = 0$  (the non-participation equilibrium),  $0 < \pi < 1$  (an interior equilibrium), or  $\pi = 1$  (a boundary equilibrium). Combining equations (2)–(4), this locus is determined by:

$$V_C = (\text{or } <) V_M \Leftrightarrow (\beta - \pi a)/[\pi(b - \pi c)] = (\text{or } <) g(\beta). \tag{6}$$

This is a continuous curve which includes the points  $(x_{\min}, x_{\min}/a)$  and  $(x_{\max}, 0)$ . It divides the  $(\beta, \pi)$  space into three regions: the set of points *on* the locus, at which  $V_C = V_M$ ; the set of points *inside* the locus, at which  $V_C > V_M$ ; and the set of points *outside* the locus, at which  $V_M > V_C$ .

Propositions (5) and (6) together imply the following result about the relative positions of the two loci:

$$\text{if } E_D = E_N \text{ and } \beta < x_{\max}, \text{ then } E_C = (\text{or } <) E_M \Leftrightarrow \beta = (\text{or } <) ab/c. \quad (7)$$

The loci intersect at the non-participation equilibrium  $(x_{\max}, 0)$ . If  $x_{\max} \leq ab/c$ , there is no other intersection and hence no other equilibrium. This case is illustrated in Figure 1a. (The loci are shown by the curves ND and CM; N is the non-participation equilibrium. The arrows refer to the dynamic analysis, which will be presented in Section 4.) If instead  $x_{\max} > ab/c$ , there is exactly one other intersection, at  $\beta = ab/c$ . There are now three alternative cases.

In the first case, illustrated in Figure 1b, this intersection is at  $\pi < 1$ . This intersection, denoted I, is an interior equilibrium, defined by  $\beta = ab/c$ ,  $\pi = a/g(ab/c)$ .<sup>3</sup> These values of  $\beta$  and  $\pi$  imply that the probability with which C is played, conditional on participation in the game (i.e. conditional on N *not* being played) is  $c/(a + c)$ , ensuring that  $V_D = 0$ . (Equivalently, the frequencies with which C and D are played are in the ratio  $c:a$ .) In this case, there may also be boundary equilibria; these occur if the CM locus intersects the line  $\pi = 1$  to the left of the ND locus.

In the second case, the loci intersect at  $\pi > 1$ . Because the CM locus is continuous, and because  $x_{\min}/a < 1$ , there must be at least one point to the left of the ND locus at which the CM locus intersects the line  $\pi = 1$ . Any such point is a boundary equilibrium. This case is illustrated in Figure 1c; B is a boundary equilibrium. In the third case (not illustrated), the loci intersect exactly at  $\pi = 1$ . Then this intersection is a boundary equilibrium. In this case, there may be other boundary equilibria.

---

<sup>3</sup> The equilibrium value of  $\pi$  can be derived from (2) by using the fact that, in an interior equilibrium,  $V_D = 0$ .

The foregoing argument establishes:

*Result 1.* If  $x_{\max} > ab/c$ , there is at least one (interior or boundary) equilibrium with  $0 < \pi \leq 1$  and  $x_{\min} < \beta < x_{\max}$ .

In other words, provided the upper tail of the distribution of cooperative benefit is not too short, there is at least one equilibrium in which both C and D are played with positive probability.

We now consider players' payoffs in such equilibria. Let  $V^*(\beta, \pi)$  be the ex ante expected payoff to any player  $i$ , prior to the realisations of random variable  $X$ , given that  $i$  and his opponent play according to  $\beta$  and  $\pi$ . We will call  $V^*(\beta, \pi)$  the *value* of the game conditional on  $(\beta, \pi)$ .

The following results are derived in the Appendix:

*Result 2.* In every interior and boundary equilibrium, the value of the game is strictly positive.

*Result 3.* Suppose there are two equilibria,  $(\beta, \pi)$ ,  $(\beta', \pi')$ , such that  $\beta < \beta'$ . Then  $V^*(\beta, \pi) > V^*(\beta', \pi')$ .

Result 2 establishes that in every interior and boundary equilibrium, cooperative activity creates positive net benefits relative to the benchmark of non-participation, despite the presence of cheats. If there are multiple equilibria, one of these is distinguished by its having the lowest value of  $\beta$ . (Since there can be no more than one interior equilibrium, no two equilibria have the same value of  $\beta$ .) Result 3 establishes that this is the equilibrium at which the value of the game is greatest. We will call this the *highest-value equilibrium*.

### **3. The model: comparative statics**

The frequency of cooperative behaviour that can be sustained in equilibrium depends on the distribution of cooperative benefit  $X$ . To keep the exposition simple, we analyse the effect of a

*rightward* shift from one distribution  $F$  to an unambiguously superior distribution  $G$  when there is no change in the support  $[x_{\min}, x_{\max}]$ . That is, for all  $x_{\min} < z < x_{\max}$ ,  $G(z) < F(z)$ . The values of all other parameters are held constant.

Using (5) it can be shown that if some point  $(\beta, \pi)$  is on the ND locus for the distribution  $F$ , it is inside the corresponding locus for  $G$ . Similarly, using (6), if some point  $(\beta, \pi)$  is on the CM locus for the distribution  $F$ , it is inside the corresponding locus for  $G$ . Thus, an improvement in the distribution of cooperative benefit moves both loci outwards. Figure 2 illustrates the effects of a shift in the distribution from  $F$  (inducing the loci  $\text{ND}(F)$  and  $\text{CM}(F)$ ) to  $G$  (inducing the loci  $\text{ND}(G)$  and  $\text{CM}(G)$ ).

As this diagram shows, if the game has interior equilibria for both distributions, those equilibria have the same value of  $\beta$ , namely  $ab/c$ , but the  $G$  equilibrium has a higher value of  $\pi$ . Since  $G(ab/c) < F(ab/c)$ , and since the frequencies with which C and D are played are in the fixed ratio  $c : a$ , both C and D are played with higher frequency in the  $G$  equilibrium than in the  $F$  equilibrium. More intuitively, the relationship between cooperation and cheating is analogous to that between prey and predator. If the distribution of cooperative benefit becomes more favourable, a higher frequency of cooperation is induced; but the more cooperation there is, the more cheating can be sustained.

If the game has boundary equilibria for both distributions, the highest-value  $G$  equilibrium must be to the left of the highest-value  $F$  equilibrium. (This can be seen by considering the effect of an outward shift of the CM locus in Figure 1c.) Thus, the former equilibrium induces a higher frequency of cooperation than the latter.

The following general result is proved in the Appendix:

*Result 4.* Suppose  $x_{\max} > ab/c$  and let  $F, G$  be two distributions of  $X$  such that  $G$  is rightward of  $F$ . Then in the highest-value  $G$  equilibrium, the frequency of cooperation and the value of the game are both strictly greater than in the highest-value  $F$  equilibrium.

Thus, as the distribution of cooperative benefit becomes progressively more favourable, the maximum sustainable frequency of cooperation increases.<sup>4</sup> Increases in cooperation are associated with increases in cheating until the frequency of non-participation falls to zero.

#### 4. The model: dynamics

We now consider the dynamics of the model, under the simple assumption that  $\beta$  and  $\pi$  evolve independently. (In a biological application, this is equivalent to assuming that  $\beta$  and  $\pi$  are determined by distinct genes.) It is sufficient to assume that, in the population as a whole, the value of  $\beta$  tends to increase (respectively: decrease) if  $V_M > V_C$  ( $V_M < V_C$ ), and that the value of  $\pi$  tends to increase (decrease) if  $V_D > V_N$  ( $V_D < V_N$ ). This gives the dynamics shown in phase-diagram form in Figure 1.

We begin by considering evolutionary stability. It is immediately obvious that the non-participation equilibrium is *not* evolutionarily stable: for example, it can be invaded by any strategy that sometimes cooperates and never defects. In contrast, all interior and boundary equilibria are evolutionarily stable. In any such equilibrium, the value of  $\beta$  is *uniquely* optimal for each player, given the behaviour of the others. Thus, in analysing evolutionary stability, it is sufficient to consider the vertical arrows in the phase diagrams. It is easy to see that all interior and boundary equilibria are stable with respect to vertical movements.

Considering the dynamics more explicitly, Figures 1a, 1b and 1c all show that the non-participation equilibrium N is not locally stable. There are evolutionary paths leading to this equilibrium (from outside the CM locus) but also paths leading away from it (from inside that locus).

Figure 1b shows that in the neighbourhood of an interior equilibrium (I), the dynamics exhibit cycles. Described in terms of the evolution of the frequencies of the three strategies N, C and D, these cycles are similar to those of the Rock–Scissors–Paper game. (The frequency

---

<sup>4</sup> This comparative-static property is compatible with evidence that in both human and non-human interaction, the level of cooperation is greater, the higher the payoffs to cooperation (Heinrich et al., 2001; Clutton-Brock, 2002).

of cooperation is greatest towards the left of the diagram, where the value of  $\beta$  is low. From there, evolutionary paths lead towards the top right, where the values of  $\beta$  and  $\pi$  are both high, and the frequency of cheating is greatest. From there, paths lead towards the bottom right, where  $\beta$  is high and  $\pi$  is low, and the frequency of non-participation is greatest. And from there, paths lead back towards the left.) These paths resemble predator–prey cycles, cheats acting as predators and cooperators as prey.

If the CM locus cuts the line  $\pi = 1$  at a point where  $\beta < ab/c$ , this point is a boundary equilibrium. If (as in the case shown in Figure 1c) points to the left of this equilibrium are outside the locus, the equilibrium is locally stable. Not all boundary equilibria have this property, but whenever the ND and CM loci intersect at  $\pi > 1$ , there must be at least one locally stable boundary equilibrium.

## 5. Discussion

We do not intend to claim that our model represents *the* mechanism that underlies human and animal cooperation. There is no good reason to suppose that cooperation is a single phenomenon with a unified causal explanation. We find it more plausible to view cooperation as a family of loosely-related phenomena which may have multiple causes. We offer our model as a stylised representation of *one* mechanism by which cooperation might emerge and persist.

Our model is unusually robust in that it assumes only materially self-interested motivations and applies to anonymous, well-mixed populations. In claiming this as a merit of the model, we do not deny that individuals sometimes act on pro-social motivations, nor that many recurrent cooperative interactions are between individuals who are known to one another, nor that populations of potential cooperators are often structured into clusters of individuals who interact mainly with their ‘neighbours’. Each of these factors can contribute to the explanation of cooperation in particular environments. Nevertheless, theories that depend on non-anonymity, or on population structures taking particular forms, have restricted domains of application. And it is hardly controversial to claim that self-interest is a particularly common and reliable motivation.

As an illustration of how theories with less robust assumptions can be restricted in their application, we consider the currently much-discussed hypothesis of altruistic punishment (Fehr and Gächter, 2000; Gintis et al, 2005). The key insight is that multilateral cooperation can be sustained in equilibrium if individuals have low-cost options of punishing one another, and if even a relatively small proportion of individuals have relatively weak preferences for punishing non-cooperators. However, the general effectiveness of this mechanism depends on the cost of punishing being low relative to the harm inflicted, and on the absence of opportunities for punishees to retaliate (Herrmann et al, 2008; Nikiforakis, 2008); and it requires that at least some individuals have non-selfish preferences for punishing. Such preferences might be sustained by *cultural* group selection in hunter-gatherer economies, where groups are small and inter-group warfare is frequent, but these conditions are not typical of the modern world; even among hunter-gatherers, *biological* group selection of altruistic punishment would be frustrated by inter-group gene flow (Boyd et al, 2005). Altruistic punishment should be understood as a mechanism that can sustain cooperation in specific types of environment, not as *the* solution to the problem of explaining cooperation. We claim no more than this for our own model.

We have said that our model is in the same spirit as some recent work by biologists, which finds apparently cooperative behaviour to be directly beneficial to the individual cooperator (see Section 1 above). But, as we now explain, the explanatory principles used by these biologists are not the same as those exhibited in our model.

One of the fundamental features of our model is that the cooperative behaviour it describes is *reciprocally beneficial*. By this, we mean the following. Such cooperation is not simply a unilateral action by one individual which, intentionally or unintentionally, confers benefits on another; it is the *composition* of cooperative actions by two or more individuals, the combined effect of which is to benefit each of them. In other words, each cooperator benefits from his action *only if* this action is reciprocated by one or more other individuals. In the absence of enforceable promises, reciprocally beneficial cooperation requires at least one individual to choose a cooperative action without assurance that others will reciprocate. In our model, any player who chooses to cooperate incurs a risk of loss, which is realised if his opponent cheats. One might think (as we are inclined to do) that reciprocal benefit is a

hallmark of genuine, as opposed to apparent, cooperation (see also Sachs et al., 2004; West et al. 2007). In biological models of mutualism, cooperation is not reciprocally beneficial, in the sense we have defined.

In the Snowdrift game, which is often used to model apparently cooperative animal behaviour, cooperation and cheating are best responses to one another. In the original story, two drivers are stuck in the same snowdrift. Both drivers have shovels, and so each can choose whether or not to dig. If either driver digs a way out for his own car, the other can drive out too. Each would rather be the only one to dig than remain stuck. This defines a game with Chicken payoffs; in a pure-strategy Nash equilibrium, one driver digs and the other free-rides (Sugden, 1986). Such an equilibrium is not a case of reciprocally beneficial behaviour.

Clutton-Brock (2009) offers the Soldier's Dilemma as a model of mutualism in biology. In this game, a patrol of soldiers is ambushed by the enemy. Soldiers who fire back attract incoming fire and increase their chance of being killed. By firing back, however, each individual reduces the probability that the patrol will be overrun. The gain from this may be such that from an individual's perspective there is no dilemma at all: firing back may give the best chance of individual survival, irrespective of what the others do. A biological equivalent to this game (or perhaps to Snowdrift) can be found in the behaviour of certain birds and mammals, such as Arabian babblers and meerkat, which feed in predator-rich environments. Individuals of these species go on sentinel duty once they have fed for long enough to be close to satiation (Clutton Brock et al., 1999). In these games, cooperation is chosen either as a dominant strategy or as a best response to other players' non-cooperation; it is not reciprocally beneficial.

In the story of the Soldier's Dilemma, it would be natural to assume that cooperation would be a dominant strategy only if the number of soldiers in the patrol was relatively small, so that each of them received a significant share of the total benefit created by his own cooperative action. Hauert et al (2002) present a model which can be understood as a version of the Soldier's Dilemma in which the size of the patrol is endogenous. This is an  $n$ -player model of voluntary contributions to a public good, but with an outside option of non-participation. A player who takes the outside option receives a small positive payoff  $\sigma$  with



certainty, but forgoes any share in the benefits of the public good. Players who participate can either cooperate (contribute to the public good) or cheat (not contribute). Each cooperator incurs a cost of 1 and creates a benefit of  $r$  (where  $1 < r < n$  and  $r > \sigma + 1$ ), which is divided equally between all participants. This game has no pure-strategy Nash equilibrium. (If all one's opponents take the outside option, the best response is to cooperate; if they all cooperate, the best response is to cheat; if they all cheat, the best response is the outside option.) There is a unique symmetrical mixed-strategy Nash equilibrium in which the expected payoff to all three strategies is  $\sigma$ . More intuitively, in equilibrium the expected number of participants in each game is sufficiently small that cooperation and cheating are equally profitable. Replicator dynamics have the Rock-Scissors-Paper cyclical pattern.

There are some similarities between Hauert et al's model and ours: both models include a non-participation option, and both induce mixed-strategy equilibria with predator-prey characteristics. However, Hauert et al's model differs from ours in two significant ways. First, the mechanism that induces cooperation works through variation in the number of participants in the cooperative activity. For this reason, the model cannot represent cooperative activities which require a fixed number of participants. In particular, it cannot represent activities which inherently involve just two individuals – as, for example, most forms of market exchange do. Second, because the costs and benefits of contributing to the public good are non-stochastic, the expected payoffs to cooperation, cheating and non-participation are equal in equilibrium. Thus, although some cooperative activity takes place in equilibrium, this activity generates no net benefit relative to non-participation: it is not reciprocally beneficial.

We suggest that our analysis provides a stylized but essentially realistic account of a mechanism by which reciprocally beneficial cooperation can emerge and persist in anonymous, well-mixed populations in which strategies are selected according to their material payoffs. Using two simple components – voluntary participation and stochastic payoffs – that have not previously been put together, we have constructed a robust general-purpose model of cooperation.

We are conscious that, for some theoretically-oriented economists, the mechanism we have described may seem rather prosaic. For decades, the Prisoner's Dilemma has been used

as the paradigm model of cooperation problems, and the problem of explaining cooperation in that game has been treated as a supreme theoretical challenge. Viewed in that perspective, a modelling strategy which relaxes the assumption that cooperation is always a dominated strategy may seem too easy. But we share the view of Worden and Levin (2007) that many real-world cooperation problems are less intractable than the Prisoner's Dilemma. Neglecting these cases results in an incomplete body of theory and fosters unwarranted pessimism about the possibility of spontaneous cooperation.

## Appendix: Proofs of results

*Proof of Result 2:* Let  $(\beta, \pi)$  be any interior or boundary equilibrium, and consider any player  $i$ . With probability  $F(\beta)$ ,  $x_i < \beta$  and  $i$  plays N or D. In an interior equilibrium,  $V_D = V_N = 0$ . In a boundary equilibrium,  $V_D \geq V_N = 0$  and N is not played. In either case,  $i$ 's expected payoff is equal to  $V_D$  and is non-negative. With probability  $1 - F(\beta)$ ,  $x_i \geq \beta$  and  $i$  plays C. If  $x_i = \beta$ ,  $i$  is indifferent between C and D and the expected payoff is again  $V_D$ . If  $x_i > \beta$ ,  $i$  plays C; his expected payoff (conditional on  $x_i > \beta$ ) exceeds that in the  $x_i = \beta$  case by  $[1 - F(\beta)](x_i - \beta)$ ; here  $1 - F(\beta)$  represents the probability that  $i$ 's opponent plays C. Hence:

$$V^*(\beta, \pi) = V_D + [1 - F(\beta)] E[\max(x_i - \beta, 0)], \quad (\text{A1})$$

where E is the expectation operator. Since  $V_D \geq 0$  and  $\beta < x_{\max}$ , the value of  $V^*(\beta, \pi)$  is strictly positive.

*Proof of Result 3:* If  $(\beta, \pi)$  and  $(\beta', \pi')$  are both interior and/or boundary equilibria, Result 3 can be derived from (A1) using the fact that  $V_D$  is decreasing in  $\beta$  (an implication of (2)). If  $(\beta, \pi)$  is the non-participation equilibrium,  $V^*(\beta, \pi) = 0$  and so Result 3 follows trivially from Result 2.

*Proof of Result 4:* Suppose  $x_{\max} > ab/c$ . Let  $(\beta, \pi)$  be the highest-value  $F$  equilibrium and let  $(\beta', \pi')$  be the highest-value  $G$  equilibrium. There are three possibilities. *Case 1:*  $(\beta, \pi)$  and  $(\beta', \pi')$  are both interior equilibria. Then  $\beta' = \beta = ab/c$  and  $\pi' > \pi$ . (This case is illustrated in Figure 2) Since  $G(\beta') < F(\beta)$ , the frequency of cooperation is higher in the  $G$  equilibrium. Using (A1) and the fact that  $V_D = 0$  in every interior equilibrium, it can be shown that the value of the game is strictly greater in the  $G$  equilibrium. *Case 2:*  $(\beta, \pi)$  and  $(\beta', \pi')$  are both boundary equilibria. Then (because the CM locus for  $G$  lies outside the CM locus for  $F$ )  $\beta' < \beta$  and  $\pi' = \pi = 1$ . Since  $G(\beta') < F(\beta)$ , the frequency of cooperation is higher in the  $G$  equilibrium. Using (2), it can be shown that  $V_D$  is strictly greater in the  $G$  equilibrium. Then, using (A1) in relation to the distributions  $F$  and  $G$ , it can be shown that the value of the game is strictly greater in the  $G$  equilibrium. *Case 3:*  $(\beta, \pi)$  is an interior equilibrium and  $(\beta', \pi')$  is a boundary equilibrium. Then  $\beta' \leq \beta$  and  $1 = \pi' > \pi$ . Since  $G(\beta') < F(\beta)$ , the frequency of

cooperation is higher in the  $G$  equilibrium. In the interior equilibrium,  $V_D = 0$ . In the boundary equilibrium,  $V_D \geq 0$ . Then, using (A1), it can be shown that the value of the game is strictly greater in the  $G$  equilibrium.

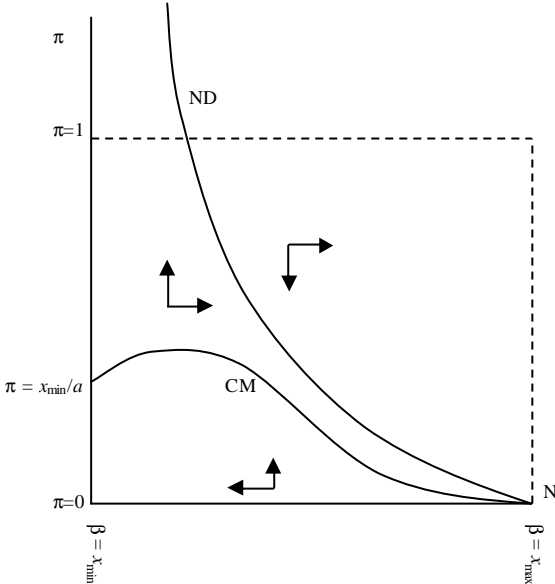
## References

- Binmore, K., 1994. *Game Theory and the Social Contract, Volume 1: Playing Fair*. MIT Press, Cambridge, MA.
- Binmore, K., 1998. *Game Theory and the Social Contract, Volume 2: Just Playing*. MIT Press, Cambridge, MA.
- Boyd, R., Gintis, H., Bowles, S. and Richerson, P.J., 2005. The evolution of altruistic punishment. In Gintis et al, 2005, pp. 215-227.
- Clutton-Brock, T.H., 2002. Breeding together: kin selection and mutualism in cooperative vertebrates. *Science* 296, 69-72.
- Clutton-Brock, T. H., 2009. Cooperation between non-kin in animal societies. *Nature* 462, 51-57.
- Clutton-Brock, T. H., O'Rian, M.J., Brotherton, P.N.M., Gaynor, D., Kansky, R., Griffin, A.S., Manser, M., 1999. Selfish sentinels in cooperative mammals. *Science* 284, 1640-1644.
- Darwin, C., 1859. *On the origin of species by means of natural selection*. John Murray, London.
- Fehr, E. and Gächter, S., 2000. Cooperation and Punishment. *American Economic Review* 90(4), 980-994.
- Gintis, H., Bowles, S., Boyd, R., Fehr, E. (eds), 2005. *Moral Sentiments and Material Interests*. MIT Press, Cambridge, MA.
- Hauert, C., De Monte, S., Hofbauer, J., Sigmund, K., 2002. Volunteering as Red Queen Mechanism for Cooperation in Public Goods Game. *Science* 296, 1129-1132.
- Henrich, J., Boyd, R., Bowles, S., Camerer, C., Fehr, E., Gintis, H., McElreath, R., 2001. In search of homo economicus: behavioural experiments in 15 small-scale societies. *The American Economic Review* 91 (2), 73-78.
- Herrmann, B., Thöni, C., Gächter, 2008. Antisocial punishment across societies. *Science* 319, 1362-1367.

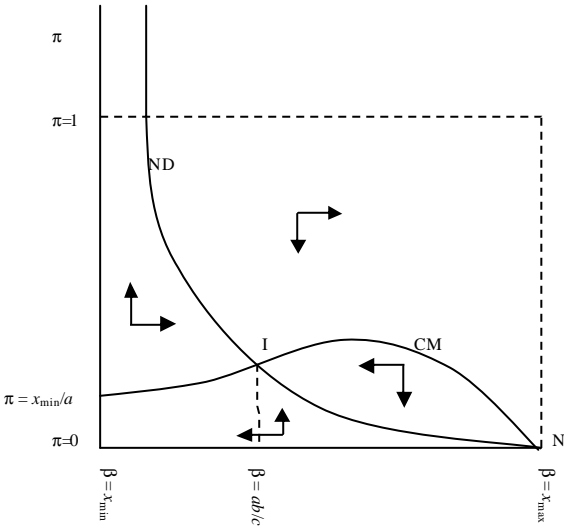
- Leimar, O, Hammerstein, P., 2001. Evolution of cooperation through indirect reciprocity. *Proceedings of the Royal Society* 268, 745-753.
- Nikiforakis, N., 2008. Punishment and counter-punishment in public good games: can we really govern ourselves? *Journal of Public Economics* 92, 91-112.
- Nowak, M.A., 2006. Five rules for the evolution of cooperation. *Science* 314, 1560-1563.
- Nowak, M.A., Sigmund, K., 2005. Evolution of indirect reciprocity. *Nature* 437, 1291-1298.
- Rousseau, J.-J., 1755/ 1998. Discourse on the origin and foundations of inequality among men. In Ritter, A., Bondanella, J.C. (eds), 1998, *Rousseau's Political Writings*. Norton, New York, pp. 3-57.
- Sachs, J.L., Mueller, U.G., Wilcox, T.P., Bull, J.J., 2004. The evolution of cooperation. *Quarterly Review of Biology* 79 (2), 135-160.
- Smith, A., 1763; 1978. *Lectures on Jurisprudence*. Oxford University Press, Oxford.
- Sugden, R., 1986. *The Economics of Rights, Cooperation and Welfare*. Blackwell, Oxford. (Second edition 2004, Palgrave-Macmillan, Basingstoke.
- Worden, L., Levin, S.A., 2007. Evolutionary escape from the prisoner's dilemma. *Journal of Theoretical Biology* 245, 411-422.

**Figure 1. Equilibria and dynamics**

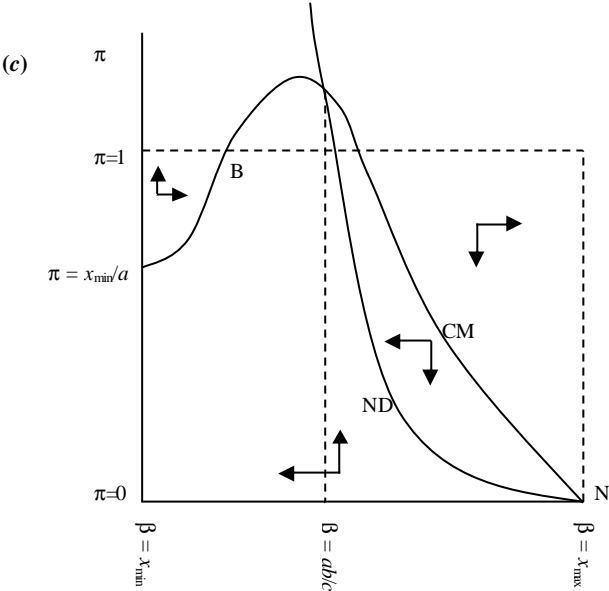
a: Non-participation the only equilibrium



b: An internal equilibrium



c: A boundary equilibrium





**Figure 2: Effects of a shift in the distribution of cooperative benefit**

