

## Emergence of Fairness in Repeated Group Interactions

S. Van Segbroeck,<sup>1</sup> J. M. Pacheco,<sup>2,3</sup> T. Lenaerts,<sup>1,4</sup> and F. C. Santos<sup>5,3</sup>

<sup>1</sup>MLG, Université Libre de Bruxelles, Brussels, Belgium

<sup>2</sup>Departamento de Matemática e Aplicações, Universidade do Minho, Braga, Portugal

<sup>3</sup>ATP-group, CMAF, Instituto para a Investigação Interdisciplinar, Lisboa, Portugal

<sup>4</sup>AI-lab, Vrije Universiteit Brussel, Brussels, Belgium

<sup>5</sup>DEI, & INESC-ID, Instituto Superior Técnico, TU Lisbon, Lisboa, Portugal

(Received 26 August 2011; published 10 April 2012)

Often groups need to meet repeatedly before a decision is reached. Hence, most individual decisions will be contingent on decisions taken previously by others. In particular, the decision to cooperate or not will depend on one's own assessment of what constitutes a fair group outcome. Making use of a repeated  $N$ -person prisoner's dilemma, we show that reciprocation towards groups opens a window of opportunity for cooperation to thrive, leading populations to engage in dynamics involving both coordination and coexistence, and characterized by cycles of cooperation and defection. Furthermore, we show that this process leads to the emergence of fairness, whose level will depend on the dilemma at stake.

DOI: 10.1103/PhysRevLett.108.158104

PACS numbers: 87.23.Kg, 89.75.Fb

Many problems of cooperation among humans boil down to the dilemma of helping others at a cost to ourselves or refraining from doing so while still profiting from the help provided by others [1–3]. Surprisingly often we take the first option, even though rational considerations encourage us not to [1,2]. This talent for cooperation forms one of the cornerstones of human society and is, as such, also largely responsible for the unprecedented success of our species [4]. But how did evolution succeed in shaping such cooperative beings, if the temptation to free ride on the benefits produced by others is always lurking? This paradox of cooperation [5] has been under intense scrutiny for decades and, fortunately, several mechanisms discourage us from actually giving in to this temptation [5–15]. Physicists have investigated some of these mechanisms (for an excellent review, see [8]), as human cooperation constitutes an excellent example of a complex system. Cooperation may, for instance, be worthwhile if your opponent has the chance to return you the favor later on. If he or she is not willing to do so, his or her cheating behavior can still be retaliated. This is Robert Trivers' *direct reciprocity* at work [16]. Theoretical and empirical studies show that individuals who pursue long-term relationships built on mutual cooperation are expected to prevail [17–21]. In this context, tit-for-tat players constitute the most famous example [17]: They always start by cooperating, subsequently repeating their opponent's last move.

Direct reciprocation may enhance cooperation for pairwise interactions, but when larger groups of actors are involved, decision-making becomes much more complex. Similar to the relation between 2-body and many-body interactions in Physics, also in human decisions there is a significant increase in complexity when going from pairwise cooperative game interactions to collective efforts in sizable groups. Technically, such an increase in complexity

is reflected in the number of possible behavioral equilibria, which scales linearly with the group size [22], even in the absence of reactive players. Moreover, it is far from clear under which conditions a cooperator (defector) should switch to defection (cooperation) when engaged in a repeated collective endeavor, wherein some may cooperate while others defect. To whom should one reciprocate [23]? One possibility is to reciprocate towards the entire group. As in previous studies of evolution and assessment of fair offers [24–27], reciprocating towards groups will depend on what is reckoned as a *fair* collective effort, as individuals may develop an aspiration level above which they cooperate, defecting otherwise. Such individuals constitute a  $N$ -person generalization of the 2-person reciprocators. Unsurprisingly, the spectrum of possible reciprocator strategies for group,  $N$ -person game interactions, is much larger than in the 2-person case. Some reciprocators may, for instance, be willing to cooperate only if the entire group did so in a previous encounter, whereas others may cooperate also in the presence of group members who defected.

Let us consider group decisions involving  $N$  individuals described in terms of the repeated  $N$ -person prisoner's dilemma (NPD) [28,29], in which all players have the opportunity to contribute a certain amount  $c$  ("cost") to the public good. The accumulated amount is multiplied by an investment factor  $F$  and subsequently shared equally among all group members, irrespective of their contribution. This entire process repeats itself with a probability  $w$ , resulting in an average number of  $\langle r \rangle = (1 - w)^{-1}$  rounds per group [5,30]. The outcome of the game may differ from round to round, as individuals can base their decision to contribute on the result of the previous round. We distinguish  $N$  different aspiration levels, encoded in terms of the strategies  $R_M$  ( $M \in \{1, \dots, N\}$ ).  $R_M$  players always contribute in the first round. Subsequently, they contribute

only if at least  $M$  players did contribute in the previous round. The threshold  $M$  can be regarded as their own perception of a *fair* number of contributions to the public good. In addition to these  $N$  different types of reciprocators, we include the strategy AD (always defect) to account for unconditional defectors.

Let us start by assuming an infinitely large population of individuals, where a fraction  $x$  of the population plays  $R_M$ —allowing one single value of  $M$  in the set of all reciprocators—while the remaining fraction plays AD. This will allow us to define the notation before addressing finite populations and an evolving  $M$ , while unveiling a dynamical scenario which differs strongly from the one obtained from (repeated) 2-person games. Behavioral dynamics often relies on individuals' propensity to be influenced by the actions and achievements of others. Such social learning or evolutionary dynamics can be described by the replicator equation [31]  $\dot{x} = x(1-x)(f_{R_M} - f_{AD})$ , where  $f_{R_M}$  ( $f_{AD}$ ) stand for the fitness—or success—of  $R_M$  (AD) players, given by their payoff derived from the game group interactions (see Eq. 1 in [32]).

A little algebra allows us to show that the (deterministic) replicator dynamics leads to scenarios in which cooperation may prevail, in connection with at most two internal fixed points, associated with unstable (coordination,  $x_L^*$ ) and stable (coexistence,  $x_R^*$ ) equilibria, which depend on the values of  $w$ ,  $M$ ,  $N$ , and  $F$  (for detailed derivations, see Section 1 in [32]). Intuitively, the simultaneous occurrence of these two equilibria, which happens when we face a repeated NPD ( $F < N$ ), can be explained as follows. If the  $R_M$  frequency is smaller than  $x_L^*$ , there are only a few groups in which  $R_M$  players remain cooperative for the entire duration of the game. The benefits they receive from such interactions are insufficient to cover the cost for always being prepared to cooperate in the first round, making them disadvantageous with respect to ADs. Hence,  $R_M$  players will only endure as long as their prevalence remains above a minimum fraction  $x_L^*$ , representing an unstable fixed point (coordination). But even if they succeed, they will never take over the entire population, unless  $M = N$ . As long as  $M < N$ ,  $R_M$  players will cooperate in partially cooperative groups, opening an escape hatch to the survival of a small fraction of ADs (a fraction  $1 - x_R^*$ ), reflecting the stable coexistence between the two strategies.

Additional insight in the characterization of  $x_L^*$  and  $x_R^*$  is provided in Fig. 1 which shows that, for given  $M/N$ , there is a critical probability  $\bar{w}$  above which the two equilibria emerge.  $\bar{w}$  increases as we reduce  $M/N$ , meaning that more rounds are required to prevent AD from dominating the population. Naturally, the location of  $x_L^*$  and  $x_R^*$  follows the same trend, creating an interesting, but delicate, balance between the size of the basin of attraction of the coexistence state  $x_R^*$  and its actual value. Relaxing the criterium of fairness for reciprocators (lowering  $M$ ) makes

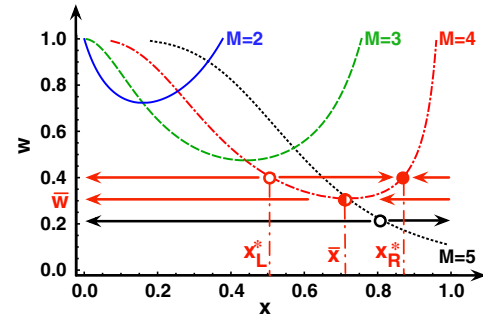


FIG. 1 (color online). Interior fixed points of the replicator equation as a function of  $w$  and  $F$ . (a) Each curve shows the position of the internal roots of the replicator equation as a function of  $w$  for a particular value of  $M$ . There are no roots if  $w$  is smaller than the critical value  $\bar{w}$ . A further increase in  $w$  leads either to two fixed points (if  $M < N$ ), the left one being unstable and the other one stable, or to just one unstable fixed point (if  $M = N$ ). The arrows indicate the direction of selection ( $F = 3.5$ ,  $N = 5$ ). Results for finite populations with evolving  $M$  are shown in Figs. 2 and 3.

the cooperative basin of attraction easier to reach (by reducing  $x_L^*$ ), but less cooperative overall (reduction of  $x_R^*$ ).

So far, we have investigated the competition between a single type of reciprocators and unconditional defectors. However, the assessment of what constitutes a fair level of cooperation in a group does not need to be unanimous in the population: The value of  $M$  itself may be under selective pressure, and in this case, the delicate competition just described becomes particularly important, mostly if we take into account that populations are finite [11] and selection is not free from errors of imitation [12,13,33] and behavioral mutations [14].

Let us then consider a population of finite size  $Z$ , and compute the average prevalence of each of the  $N + 1$  available strategies—AD plus the  $N$  different  $R_M$  strategies—over time. We implement a stochastic, finite population analogue of the deterministic evolutionary dynamics defined before, in which strategies evolve according to a mutation-selection process defined in discrete time. At each time step, the strategy of one randomly selected individual  $A$  is updated. With probability  $\mu$ ,  $A$  suffers a mutation, adopting a strategy drawn randomly from the space of  $N + 1$  available strategies; with probability  $1 - \mu$ , another randomly selected individual  $B$  acts as a role model for  $A$ : The probability that  $A$  adopts the strategy of  $B$  is given by the Fermi distribution  $p = [1 + e^{\beta(f_A - f_B)}]^{-1}$  [8,12,33,34], where  $f_A$  ( $f_B$ ) denotes the fitness of individual  $A$  ( $B$ ) and  $\beta \geq 0$  measures the strength of the fitness contribution to the update process, i.e., the so-called *intensity of selection* [12].

In the limit in which mutations are rare, we are able to compute analytically the relative prevalence of each of the different strategies [15,21,35] (for details, see Section 2 of [32]). This simplified limit turns out to be valid over a

much wider interval of mutation regimes, as we show below via numerical simulations. In this limit, the population will either end up wiping out the mutant or witness the fixation of the intruder long before the occurrence of a new mutation. Hence, there will be a maximum of two strategies present simultaneously in the population. The fixation probabilities of all possible mutants in all (otherwise) monomorphic populations can be readily computed analytically [12,15,32], defining a reduced (embedded) Markov chain, with which we compute the stationary distribution of the population, i.e., the average fraction of time the population spends in each of the  $N + 1$  monomorphic configurations of the population [32,35].

The results are shown in Fig. 2(a), where we plot the stationary distribution for different values of the parameter  $F$ . The distribution of  $R_M$  players reveals some remarkable features: On one hand, there is a specific concept of fairness, associated with an aspiration level  $M^*$  whose corresponding strategy  $R_{M^*}$  is most favored by evolution, being the most prevalent among all  $R_M$  strategies. On the other hand, unanimity in the assessment of fairness does not occur, given that several values of  $M$  may coexist in the population. Finally, as the dilemma becomes harsher

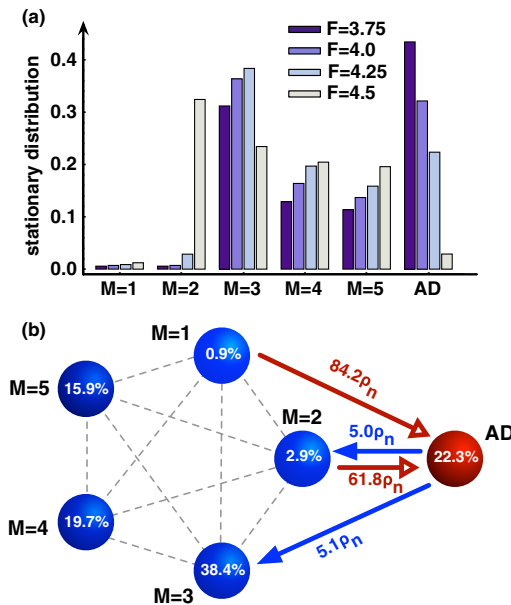


FIG. 2 (color online). Evolution of the threshold  $M$  in a finite population. (a) Stationary distribution for different values of  $F$ . Irrespective of the success of AD players, which decreases with increasing  $F$ , there is always an optimal threshold  $M^*$ , whose corresponding strategy  $R_{M^*}$  is the most prevalent ( $w = 0.9$ ,  $N = 5$ ,  $Z = 100$ ,  $\beta = 1.0$ ). (b) The percentages indicate the fraction of time the population spends in each composition of the population ( $F = 4.25$ ). Arrows indicate transitions whose fixation probability is greater than  $\rho_N = 1/Z$ . One observes oscillations between cooperation and defection. The population moves from  $R_M$  with small  $M$ , over AD, back to  $R_M$  with moderate threshold. Neutral drift may bring us back to  $R_M$  with small  $M$ , as emphasized using grey dotted lines.

(lower values of  $F$ ), the higher the fraction of the population that adopts the most prevalent assessment of fairness— $M^*$ —which is always much smaller than the group size  $N$  [see also Fig. 3(a)]. Naturally, the success of AD players increases with decreasing  $F$  [36].

The intuition behind the emergence of an optimal level of fairness  $M$  can be understood with the help of Fig. 2(b), where we analyze the typical flow of probability between the different monomorphic states. Arrows represent transitions favored by natural selection, i.e., those whose fixation probability exceeds  $1/Z$  (associated with the fixation probability of a mutant under neutral evolution). Suppose we start from a homogeneous population of ADs. Figure 2(b) shows that there are several  $R_M$  types with intermediate  $M$  who can invade AD (solid blue arrows). Clearly, such a modest assessment of what constitutes a fair group (intermediate  $M$ ) combines the best of two worlds: avoiding continuous exploitation, but being sufficiently generous to maintain the level of cooperation in groups that are only partially cooperative. Once a given  $R_M$  takes over, neutral drift (grey dashed lines) can drive the population to any of the  $N - 1$  other  $R_M$  states, which provides a foundation for the coexistence of different concepts of fairness in the population. Whenever the demands for fairness are too modest ( $M < M^* = 3$ ), AD can take

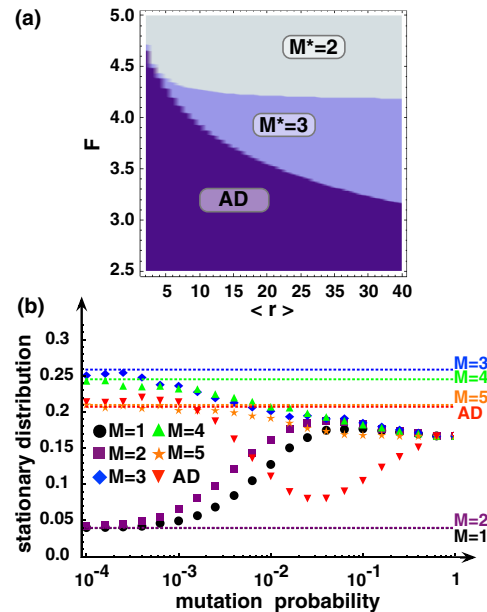


FIG. 3 (color online). Evolutionary dynamics for a) arbitrary number of rounds and b) mutation probabilities. (a) The optimal threshold  $M^*$  as a function of  $F$  and  $\langle r \rangle$  in the limit of rare mutations ( $w = 0.9$ ,  $N = 5$ ,  $Z = 100$ ,  $\beta = 1.0$ ). (b) Dashed lines indicate the stationary distribution in the small-mutation limit. Each symbol indicates, for a given mutation probability, the fraction of the population that adopts the corresponding strategy, averaged over the simulation time (30 simulations, each lasting for  $10^9$  iterations;  $w = 0.9$ ,  $N = 5$ ,  $Z = 100$ ,  $\beta = 0.05$ ; our results are robust to changes in  $\beta$  [32]).

over the population again (open red arrows). Hence, the population oscillates continuously between cooperation and defection, resembling the cycles of war and peace similar in spirit to those identified in the context of repeated 2-person games of cooperation [21].

This scenario constitutes a general feature of the present model, and is not the result of a particular choice of the average number of rounds  $\langle r \rangle$  (or  $w$ ) or mutation probability  $\mu$ , as demonstrated in Fig. 3. Figure 3(a) shows that, irrespective of the number of rounds  $\langle r \rangle$ , AD abounds when  $F$  is small,  $R_M$  with  $M = 2$  when  $F$  is large (but still smaller than the group size  $N$ ), and  $R_M$  with  $M = 3$  for intermediate values of  $F$ , which corresponds exactly to the findings reported in Fig. 2(a). In other words, evolution shapes the population assessment of fairness, depending on the constraints imposed by the collective dilemma. In Fig. 3(b) we investigate the robustness of our results with respect to changes in  $\mu$ . We abandon the limit of rare mutations, and determine the stationary distributions for arbitrary mutation rates via computer simulations. For  $\mu < Z^{-2}$ , the results match the limit of rare mutations. More importantly, the plot shows that our general conclusion remains valid for a wide range of mutation probabilities:  $R_M$  players with a moderately large aspiration are expected to prevail throughout a wide range of mutation values. For large mutation rates ( $\mu > Z^{-1}$ ), all types of reciprocators become equally probable and dominant with respect to ADs. As a result, the overall outcome of cooperation is enhanced for high mutation rates. This is an important point, as one expects that, e.g. in human interactions, errors of decision making, well captured by the behavioral mutations introduced here, may be sizable [14], although at present a quantitative estimate is lacking. Needless to say, the results shown in Figs. 2 and 3 for  $N = 5$ , remain valid for other values of  $N$ , in the sense that the physical order parameter of the model remains the ratio  $M/N$  (see  $\bar{x}$  in Section I of [32]).

In summary, we have studied the evolutionary dynamics of repeated group interactions, in which individuals engage in an iterated NPD. Reciprocators are defined as individuals who may cooperate, contingent on their own individual assessment of what constitutes a fair group contribution. We found that evolution selects for a moderate, yet prevalent, concept of fairness in the population. This choice results from a detailed competition between the capacity to avoid continuous exploitation and the generosity of contributing in groups which are only partially cooperative. The prevalent concept of fairness that emerges in the population constitutes a compromise between too low aspiration levels, which lead reciprocators to extinction, and too high aspiration levels, associated with harsh coordination thresholds. Combined with the neutrality between different concepts of fairness, the emergent dynamics leads to cyclic behavior which, being ubiquitous in evolutionary games [8,33], also resembles the alternation

between cooperation and defection which seems to pervade throughout human history [37].

Financial support from FNRS Belgium (S.V.S., T.L.) and FCT-Portugal (F.C.S., J.M.P.) is gratefully acknowledged.

- 
- [1] G. Hardin, *Science* **162**, 1243 (1968).
  - [2] P. Kollock, *Annu. Rev. Sociol.* **24**, 183 (1998).
  - [3] M. Olson, *The Logic Of Collective Action: Public Goods and the Theory Of Groups* (Harvard University Press, Cambridge, MA, 1971).
  - [4] J. Maynard Smith and E. Szathmary, *The Major Transitions in Evolution* (Freeman, Oxford, 1995).
  - [5] K. Sigmund, *The Calculus of Selfishness* (Princeton University Press, Princeton, NJ, 2010).
  - [6] F.C. Santos and J.M. Pacheco, *Phys. Rev. Lett.* **95**, 098104 (2005).
  - [7] J. Gomez-Gardenes, M. Campillo, L.M. Florıa, and Y. Moreno, *Phys. Rev. Lett.* **98**, 108103 (2007).
  - [8] G. Szabo, G. Fath, *Phys. Rep.* **446**, 97 (2007).
  - [9] M. Perc and A. Szolnoki, *Phys. Rev. E* **77**, 011904 (2008).
  - [10] G. Szabo and C. Hauert, *Phys. Rev. Lett.* **89**, 118101 (2002).
  - [11] A. Traulsen, J.C. Claussen, and C. Hauert, *Phys. Rev. Lett.* **95**, 238701 (2005).
  - [12] A. Traulsen, M. A. Nowak, and J. M. Pacheco, *Phys. Rev. E* **74**, 011909 (2006).
  - [13] A. Traulsen, J. M. Pacheco, and L. A. Imhof, *Phys. Rev. E* **74**, 021905 (2006).
  - [14] A. Traulsen *et al.*, *Proc. Natl. Acad. Sci. U.S.A.* **106**, 709 (2009).
  - [15] S. Van Segbroeck, F.C. Santos, T. Lenaerts, and J.M. Pacheco, *Phys. Rev. Lett.* **102**, 058105 (2009).
  - [16] R. Trivers, *Q. Rev. Biol.* **46**, 35 (1971).
  - [17] R. Axelrod and W.D. Hamilton, *Science* **211**, 1390 (1981).
  - [18] M. Milinski, *Nature (London)* **325**, 433 (1987).
  - [19] D. Fudenberg and E. S. Maskin, *Am. Econ. Rev.* **80**, 274 (1990).
  - [20] M. A. Nowak and K. Sigmund, *Nature (London)* **355**, 250 (1992).
  - [21] L. A. Imhof, D. Fudenberg, and M. A. Nowak, *Proc. Natl. Acad. Sci. U.S.A.* **102**, 10797 (2005).
  - [22] C. S. Gokhale and A. Traulsen, *Proc. Natl. Acad. Sci. U.S.A.* **107**, 5500 (2010).
  - [23] K. Sigmund, *Trends Ecol. Evol.* **22**, 593 (2007).
  - [24] M. Nowak, K. Page, and K. Sigmund, *Science* **289**, 1773 (2000).
  - [25] J. Henrich, R. Boyd, S. Bowles, C. Camerer, H. Gintis, R. McElreath, and E. Fehr, *Am. Econ. Rev.* **91**, 73 (2001).
  - [26] K. Sigmund, E. Fehr, and M. Nowak, *Sci. Am.* **286**, 82 (2002).
  - [27] J. Henrich *et al.*, *Science* **327**, 1480 (2010).
  - [28] R. Boyd and P.J. Richerson, *J. Theor. Biol.* **132**, 337 (1988).
  - [29] S. Kurokawa and Y. Ihara, *Proc. Biol. Sci.* **276**, 1379 (2009).
  - [30] In this setting,  $w$  and  $\langle r \rangle$  can be used interchangeably.

- [31] J. Hofbauer and K. Sigmund, *Evolutionary Games and Population Dynamics* (Cambridge University Press, Cambridge, England, 1998).
- [32] See Supplemental Material at <http://link.aps.org/supplemental/10.1103/PhysRevLett.108.158104> for details.
- [33] A. Szolnoki, M. Perc, and G. Szabó, *Phys. Rev. E* **80**, 056109 (2009).
- [34] G. Szabó and C. Toke, *Phys. Rev. E* **58**, 69 (1998).
- [35] D. Fudenberg and L. Imhof, *J. Econ. Theory* **131**, 251 (2006).
- [36] In the presence of execution errors a similar result is obtained, in which the emergent  $M^*$  depends on  $F$ ,  $w$ , and the fraction of errors.
- [37] P. Turchin, *War and Peace and War: The Life Cycles of Imperial Nations* (Pi Press, New York, 2006).