# GAMES EVOLUTION PLAYS

Karl Sigmund

Institut für Mathematik
Universität Wien
Strudlhofgasse 4, A-1090 Vienna, Austria

## ABSTRACT

The evolution of cooperation is frequently analysed in terms of the repeated Prisoner's Dilemma game. Computer simulations show that the emergence of cooperation is a robust phenomenon. However, the strategy which eventually gets adopted in the population seems to depend sensitively on fine details of the modelling process, so that it becomes difficult to predict the evolutionary outcome in real populations.

## 1. EVOLUTIONARY GAME THEORY

Every science has a predilection for singularities, for phenomena which are not buried in the ceaseless stream of day-to-day routine, but stand distinctly out. Physicists addressed the motion of heavenly bodies long before they turned to falling leaves or bubbles of foam or the mechanics of human locomotion, and economists dealt with the impact of monetary devaluation upon national wealth centuries before they discovered the continuous trading of services and goods that pervades the humblest household. One reason is, of course, that a phenomenon which stands out is likely to be fairly isolated from its environment, and therefore easier to analyse. There are far fewer forces acting on a planet than on a knee-cap.

It took economists a long time to appreciate that their science has to be founded on day-to-day individual behaviour. The tendency to act in groups—as firms, or guilds, or nations, and so on—is so widespread that such groups were often viewed as natural units. They are not, of course: they result from the aggregation—spontaneous or enforced—of human individuals, and have to be understood on this level. This point took a surprisingly long time to sink in. A similar delay occured in evolutionary biology: group or species selection arguments held their own, against gene or individual selection arguments, with a considerable tenacity.

Interestingly, it was the same methodological instrument that helped, both in mathematical economics and in theoretical biology, to bolster the approach based on the level of the individual. This is the use of *game theory*, a mathematical discipline founded more

than fifty years ago by John von Neumann and Oscar Morgenstern, who were motivated by games like poker and chess to attempt an analysis of all sorts of conflicts of interest. However, note that many of the high hopes originally raised by game theory as a means of solving social, economic and military conflicts have been disappointed. But later, after game theory was amended in some essential aspects by John Nash and John Maynard Smith, it provided an extremely useful framework for discussing the social results of individual moves. Roughly speaking, John Nash introduced the main equilibrium concept to analyse non-cooperative games, which allowed to deal with non-zero sum situations; and John Maynard Smith formulated a population dynamics which permitted to get rid of the rationality assumption. We will encounter these two ideas in the following investigation of *reciprocal aid*, which is a basic concept both in economy and in ethology, and as essential to the social sciences as is chemical binding to chemistry.

But first, let us briefly describe the main aspects of the new *evolutionary* game theory (see Maynard Smith, 1982, Hofbauer and Sigmund, 1988, Weibull, 1996). In its simplest form, one no longer assumes that the players engaged in such games are endowed with perfect foresight and rationality, able to find their best moves and to out-guess the best moves of their adversaries. Rather, one considers players programmed to one specific type of inborn behaviour. Thus a strategy is no longer viewed as a sequence of well-plotted moves, but as a hard-wired program, a behavioural phenotype. Depending on their strategy, and that of their adversary, the players obtain a higher or lower payoff. This payoff, now, is not money, and it has little to do with an individual utility scale. It is simply an increment in Darwinian fitness, i.e. in average reproductive success. Depending on their payoff, the players have therefore more or less offspring, and the offspring inherit their strategies. Thus the frequencies of the strategies will change in the population. This, of course, can mean that the success of a strategy can change, from generation to generation. A strategy is said to be evolutionarily stable if, whenever it prevails in the population, a small minority playing an alternative strategy cannot invade under the influence of natural selection. But a given game need not admit any evolutionarily stable strategy at all; alternatively, it can admit several such strategies; and it is not always the case that a population close to an evolutionarily stable strategy will converge to it under the combined effects of mutation and selection. It is quite possible that the evolutionary path will not settle down to an equilibrium. Game theorists used to show this by means of a fictitious example called the Rock-Scissors-Paper game—three strategies in a cyclic rank ordering (Rock beats Scissors, Scissors beats Paper, Paper beats Rock). It was well understood that this was a mathematical oddity which ought never to worry real ethologists; but a few months ago, it was found that the side-blotched lizard engages in it (see Sinervo and Lively, 1996). More precisely, there exist three morphs—male mating strategies conveniently associated with different throat colours—which are exactly in such a cyclic ordering. Males with dark blue throats defend small territories. They can be invaded by more agressive males with orange throats. But these more aggressive males can no longer control their larger territories efficiently, and can in their turn be invaded by 'sneaker' males with yellow stripes (looking like females). Once there are many sneakers around, the more modest blue throated males can spread again, etc. Similar dynamic complexities may well be expected whenever there are lots of conceivable strategies—as, for instance, in most issues studied in human ethology.


## 2. MUTUAL AID AND THE PRISONER'S DILEMMA

Humans need tools to survive. Arguably, their most important tools are fellow humans, and an essential function of language consists in properly handling these human

tools—witness the myth of the tower of Babylon. Charles Darwin has certainly not been the first to stress the importance of cooperation in human societies when he wrote, in the *Descent of Man*, that 'the small strength and speed of man, his want of natural weapons etc are more than counterbalanced by his ... social qualities, which lead him to give and receive aid from his fellow-men.'

Mutual aid, of course, is not restricted to humans. Ethologists list striking examples of mutual aid in many species, under the headings of joint hunting, helping in fights, predator inspection, warning the flock, teaching the young, feeding and grooming. Some social insects display even higher degrees of cooperation than humans do. Many instances of mutual aid can be easily explained by kin selection—helping a relative is helping a watered down copy of oneself. Since human households are often organised along family lines, this may well account for some instances of human cooperation. But there are also households which seem to work just as well without any family ties—sailors on a fishing vessel, hunters on an expedition, robbers on the run. This kind of mutual aid has to be explained without genetics. It obviously is based on an economic exchange. But such a trading of assistance is much more vulnerable to abuse.

The first to point this out may well have been David Hume, whose *Treatise of Human Nature* ought to rank as one of the great books of human ethology. Some ten years ago, Robert Sugden (1986) wrote a most remarkable sequel, *The economics of rights, cooperation and welfare*, where he couched many of Hume's arguments in simple game theoretical terms and showed how relevant they still are for the study of human societies. Take the following passage, for instance:

'Your corn is ripe to-day, mine will be so tomorrow. 'Tis profitable for us both, that I shou'd labour with you to-day, and that you shou'd aid me tomorrow. I have no kindness for you, and know you have as little for me. I will not, therefore, take any pains upon your account; and shou'd I labour with you upon my own account, in expectation of a return, I know I shou'd be disappointed, and that I shou'd in vain depend upon your gratitude. Here then I leave you to labour alone: You treat me in the same manner. The seasons change; and both of us lose our harvest for want of mutual confidence and security.'

Let us translate this thought experiment into the language of game theory. The two farmers, now, are the players of a game. They both have two options, or strategies: to cooperate (play **C**) or not to cooperate, i.e. to defect (play **D**). Depending on their decisions are their respective payoffs. If both cooperate, they each obtain a reward $R$ which is better than the penalty $P$ obtained if both defect and lose their harvest. But if one helps the other, but receives no return, he obtains the sucker's payoff $S$ which is even lower than $P$, whereas the other player receives the highest payoff $T$, the temptation for unilateral defection: his harvest is safe, and he has spared himself the labour of helping his neighbour. In addition to this rank ordering of the payoff values $T>R>P>S$, we also have $2R>T+S$. This last condition means that the assistance yields more than it costs (if the two farmers were obliged to share their harvest, joint cooperation would have been better than unilateral exploitation).

Whenever the payoff values satisfy these two inequalities, we have what game theorists call a *Prisoner's Dilemma game* (see Trivers, 1985). It is encountered very frequently: indeed, whenever support costs less to the donor than it benefits the recipient, mutual cooperation is obviously advantageous for the two players. However, each player could do better still by not returning the other's help; and in case no such help is forthcoming, there is of course all the more reason to defect. The message, as Hume showed, is clear: no mutual aid. The Prisoner's Dilemma game epitomises the clash between what is

best from an individual's point of view and from that of a collective, a conflict that threatens countless forms of cooperation, including trade and mutual aid.

We can illustrate the two features added by Nash and by Maynard Smith at this stage already. First, the game is not zero-sum, and not even constant sum: indeed, the sum of the payoffs for the two players is larger if both cooperate than if both defect. Most parlour games are zero-sum, since the gain to one player is the loss to the other. John von Neumann, who was a passionate if somewhat amateurish poker player, elaborated his notion of a minimax solution as a central concept for zero-sum games, but this concept is of little interest for the vastly more frequent real-life interactions which are not zero-sum. Secondly, we need not assume that our players are rational and can analyse the game in advance, like Hume did, by effectively outguessing their antagonist. In the framework of evolutionary game theory, one considers (fictitious) populations consisting of programmed players—mere robots. Each of these robots is firmly wedded to one fixed strategy and will either always cooperate or always defect. They engage in a round-robin tournament of the Prisoner's Dilemma game. For each contestant, the total payoff will depend on which other players he encountered, and therefore on the composition of the population. A defector will, however, always achieve more than a collaborator would earn in its stead. We can assume that the more successful players have more offspring, which inherit their strategy (or alternatively that the less successful players have the tendency to switch their strategy and imitate their more successful rivals). The members of the newly composed population will again engage in a round-robin tournament, etc. In this caricature of an evolving population, where success means becoming more numerous, the outcome is clear: defectors will steadily increase, and eventually swamp the population. This outcome is inevitable, and has nothing to do with rationality or foresight.

## 3. RECIPROCAL ALTRUISM AND THE REPEATED PRISONER'S DILEMMA

The drab outlook of Hume's thought experiment is obviously contradicted by the many instances of cooperation found in human and animal societies. In fact, we stopped too soon with our quote of Hume. It is not true, he argues, that 'I shou'd in vain depend upon your gratitude'. In fact, I can depend on it because I can depend on your needing my services for the next harvest, or, in Hume's words, 'because I foresee that you will return my service, in expectation of another of the same kind'. But foresight is not actually needed; and neither is the 'teaching of moralists or politicians' which Hume seems to believe essential. All that is needed is the next harvest.

More precisely, as soon as one assumes that the same two players are engaged in a *repeated* Prisoner's Dilemma game, with a large random number of rounds, then the strategy of unconditional defection is no longer the inevitable outcome (see Axelrod, 1984). It can easily be seen that there exists no strategy which is best against all comers. If the opposite player, for instance, decides to always defect (or always to cooperate), it will be best to always defect. But if your adversary decides to cooperate until you defect and from then on never to cooperate again, you will be careful not to spoil the partnership: the temptation to defect in one round is more than offset by the prospect of permanently losing a cooperative partner. There is no best strategy in the repeated Prisoner's Dilemma game (see Axelrod, 1984).

One can, of course, adopt the evolutionary viewpoint again. It reduces to *drowning by numbers*: have lots of players play lots of Prisoner's Dilemma games for lots and lots of gen-

erations. More precisely, one can simulate (in a computer) populations where each individual is matched against randomly chosen opponents for a large random number of rounds of the game. Assume again that each player follows certain 'inbuilt' rules, and in this way accumulates benefits and costs, which are measured in reproductive success—the only valid currency in a Darwinian world. At the end of each generation, the players produce offspring in proportion to their overall balance. The offspring inherit their parents behaviour, except for a few mutants testing out new rules. And so on through the next generation and the next and beyond. We may watch the population evolve for as long as we care.

The main snag here is that there exist so many conceivable strategies to play the repeated game. Evolutionary trial and error can never explore all theoretical possibilities. But biological constraints help to reduce the rule-space to a manageable size. We need probably only consider strategies (a) consisting of a finite number of states (which we interpret as internal motivational variables), (b) making decisions on whether to cooperate or not in the next round, and (c) having a transition rule which leads, depending on the outcome of the current round, to the next internal state (see Nowak and Sigmund, 1995)

The famous Tit For Tat rule is an example: it starts cooperatively and then simply imitates the other player's previous move. This retaliatory rule will invade a population of defectors, as soon as its frequency has overcome some very small threshold: cooperation, then, spreads with increasing momentum through the population. A population of stern retaliators, on the other hand, cannot be invaded by defectors. But it is victimised by its own strictness: any unintentional defection (due, for example, to a mistake in implementing a move, or an erroneous perception of the other's move) entails a costly series of reprisals. Stern retaliators therefore give way to more tolerant strategies, prepared occasionally to forget or to extend an olive branch. Such generous behaviour, however, can only thrive on a back-cloth of cooperation. Since it could never spread in a defector's world, it needs Tit For Tat to pave the way. The situation is, in fact, somewhat similar to ecological succession. It leads to a cooperative community whose members are able to retaliate *judiciously*, not too much and not too little. However, such an equilibrium is not necessarily proof against an eventual increase, through random drift, of over-gentle strains of players. Once the community is sufficiently softened up, defectors can cash in and ultimately take over.

No two evolutionary chronicles played out on a computer are alike: contingency rules. But behind the vagaries of historical accident, some common traits show up (see, e.g. Lindgren, 1991, Nowak and Sigmund, 1992, 1993, 1994, Lindgren 1996, Boerlijst et al, 1996). Most conspicuous is the *bang-bang*-principle: Either almost all members cooperate almost all the time, or they almost always defect. The change from one regime to the other takes usually not more than a few hundred generations and occurs only rarely; it can be triggered by complex events—in most cases, polymorphism in the population builds up before a shift occurs—and occasionally, it can even change direction—a sharp plunge in cooperation can be reversed in mid-fall—but the population almost never settles down to a medium level of collaboration: it is all or nothing. This is a game-theoretical version of *punctuated equilibrium*: abrupt transitions between extended periods of stasis. As more and more mutant strategies are tested, the probability of being in a cooperative regime increases; the episodes of cooperation become more frequent and tend to last longer.

## 4. STABLE STRATEGIES FOR COOPERATION

There are many behavioural rules which can sustain a cooperative society. Generous Tit For Tat—always repay cooperation in kind, but tolerate defection with a certain prob-

ability—is but one example (see Nowak and Sigmund, 1992). An even more frequent out-come of the evolution towards cooperation is the strategy which has been christened, somewhat unfortunately, Pavlov (see Nowak and Sigmund, 1993). A Pavlov player starts with a cooperative move and henceforth cooperates if and only if, in the previous round, the other player used the same move than himself. This seems at first sight a strange rule. It is reasonable not to break a string of mutual cooperation, but why should one cooperate after a round where both players have defected? The rule becomes more transparent if one realises that it consists in repeating a move that has been rewarded, and in switching if the move has been punished. This rule knows how to let bygones be bygones. A mistaken de-fection between two Pavlovians leads to one round of mutual defection (since the success-ful defector persists and the sucker shifts), but in the next round, both players will make terms and resume cooperation, and continue happily with it. Intuitively, a defection is fol-lowed by a domestic quarrel and then cooperation is resumed. The Pavlov strategy embod-ies the 'win-stay, lose-shift' principle which psychologists view as the very simplest learning rule, and philosophers as the root of hedonistic morals. Natural selection has often been compared with learning by trial and error. Here is an instance where such a generation-wise 'learning process' leads a population of primitive robots to the most basic *individual* learning rule.

In a cooperative world, both Pavlovians and Generous Tit For Tat players do better than Tit For Tat, because they are immune to errors. Pavlov has an additional advantage: it ensures that the population cannot be subverted by unconditional cooperators. Indeed, af-ter a mistaken defection, Pavlov will not resume cooperation, but blithely continue to ex-ploit the sucker. This is important, since otherwise neutral drift will cause unconditional cooperators to spread and thus 'soften up' the population to such an extent that it loses its power to retaliate against exploiters.

In a world of defectors, both Pavlov and Generous Tit For Tat cannot catalyse the first step towards a cooperative society. They are too sanguine, and try too trustingly to resume cooperation. It needs stern retaliators to prepare the ground: the grim law of *An Eye For An Eye* was probably at the origin of every self-supporting community.

Do there exist strategies which are both able to invade, in small clusters, societies of defectors and able to resist, once they have reached predominance, every invasion attempt of a mutant strategy? Yes. One such strategy has been proposed by Robert Sugden. This is *Contrite Tit For Tat* (see Sugden, 1986, Boyd, 1989, and Wu and Axelrod, 1996). In con-trast to the strategies considered so far, the 'internal state' of a Contrite Tit For Tat player does not only depend on the last moves in the game, but also on the *standing* of each player, which can be good or bad. A player is in good standing if he has cooperated in the previous round, or if he has defected while provoked (i.e. while he was in good standing and the other player was not). In every other case defection leads to a bad standing. Con-trite Tit For Tat begins with a cooperative move, and cooperates except if provoked. There are three possible internal states for such a player, which we can label, in antropomorphic terms, as 'content', 'guilty' and 'provoked'.

If two Contrite Tit For Tat players engage in a repeated Prisoner's Dilemma, and if the first player defects by mistake, he loses his good standing. In the next round, he is 'guilty' and will cooperate, whereas the other player is 'provoked' and will defect without losing his good standing. Then both players will be in good standing again; both are 'con-tent' and resume their mutual cooperation in the following round.

A population of Contrite Tit For Tat players cannot be invaded—it is evolutionarily stable, at least if we assume—as we should certainly do—that there is always a small pos-sibility for making mistakes by mis-implementing a move. (This is the 'trembling hand'

approach due to the game theorist Reinhard Selten, see Selten, 1975, and Boyd, 1989.) Moreover, Contrite Tit For Tat is as good as Tit For Tat itself at invading a population of defectors—as soon as a small cluster of such players exists, it will grow with increasing momentum.

To recapitulate: evolutionary simulations of the repeated Prisoner's Dilemma show that a robust cooperative regime is highly probable. What is less clear is whether there exist theoretical or empirical grounds for favouring one specific cooperative strategy. Contrite Tit For Tat is certainly a contender. But there exist other evolutionarily stable strategies which lead to a cooperative population—for instance, Pavlov (for certain values of the payoffs $R$, $S$, $T$ and $P$), and *Remorse* (for the complementary values) (see Boerlijst et al, 1996). *Remorse* is the strategy that cooperates after a round of mutual cooperation or when in a bad standing. Statistical explorations suggest that Contrite Tit For Tat is about twice as likely to occur as an evolutionarily stable outcome than either Pavlov or Remorse. This is essentially due to the fact that Contrite Tit For Tat, which needs no catalyser to invade a society of defectors, has a head-start. (It should be noted that paradoxically, for very *low* values of the temptation payoff $T$ to defect unilaterally, the strategy *Weakling* has a chance of about one in ten to get established in the population. In this case, the average payoff per round is not the reward $R$ for mutual cooperation, but the average of $R$ and $P$, the punishment for mutual defection. *Weakling* is the strategy that cooperates only after being in bad standing.) Unfortunately, such statistical explorations seem to depend sensitively on the technical details of the numerical exploration, and cannot offer robust predictions, so far.

To give a sample of the kind of technical details, take for instance that most numerical experiments assume that the mistakes in the interactions are mis-implementations of an intended move. A bit of introspection might suggest that *mis-perceptions* of the other player's actions are at least as likely to cause trouble. As soon as we introduce this type of mistakes, Contrite Tit For Tat loses much of its lustre. If, in a match between two Contrite Tit For Tat players, one player mistakenly believes that the other is in bad standing, this leads to a sequence of mutual backbiting, just as with Tit For Tat. In contrast, Pavlov can easily be shown to be immune to errors in perception. Remorse, on the other hand, is not. Furthermore, neither Contrite Tit For Tat nor Remorse are prepared to exploit unconditional cooperators the way Pavlov does. We have seen that if neutral drift (which is usually not included in the computer simulations, but certainly very plausible in natural populations) can allow suckers to spread, the cooperative regime becomes threatened by all-out exploiters.

Another aspect which has certainly to be considered concerns the inherent plausibility of the contending strategies—their suitability for human players. We can, for instance, easily sympathize with states like 'guilty' or 'provoked'. Does that mean that the evolution of such feelings has been an outcome of selection (it need not be selection for the Prisoner's Dilemma, but possibly for another purpose: human cooperation, of course, is so important that it can certainly be expected to shape our 'internal states' to a large degree.) Many other finite state strategies for the repeated Prisoner's Dilemma are also evolutionarily stable, but their description often sounds arid and artificial. However, one could also describe Contrite Tit For Tat in such a way that it would not be recognizable as a natural strategy. Similarly, the rule for Pavlov can be either defined as 'win-stay, lose-shift', in which case it seems eminently sensible, or as 'cooperate after an $R$ or a $P$', in which case it seems silly to most people (it has, originally, been termed *Simpleton* for this reason). An argument strongly in favour of Pavlov, of course, is that the 'win-stay, lose-shift' rule can be applied in many other situations, independent of the Prisoner's Dilemma game: it could

have evolved in the context of foraging strategies, for instance. On the other hand, internal states like 'guilty' or 'provoked' do also play a role in human interactions having nothing to do with mutual help.

## 5. ALTERNATING ASSISTANCE

Another detail concerns the timing of interactions. So far, we have considered the so-called 'simultaneous Prisoner's Dilemma', where the two players make their choices in the same time instant. The reason for this is mostly conventional. If we take a closer look at Hume's example, we actually see that the two farmers cannot help each other at the same time: they have to take turns. Thus we ought to have modelled this interactions by an 'alternating Prisoner's Dilemma' (see Nowak and Sigmund, 1994, and Frean, 1994). In fact, the partners alternate in their roles of donor and recipient in many if not most instances of reciprocal help. Robert Trivers summarizes that 'reciprocal altruism is expected to evolve when two individuals associate long enough *to exchange roles frequently* as potential altruist and recipient' (Trivers, 1985).

Thus we ought also to study the iterated Prisoner's Dilemma when the players have to take turns. The slight modification in such an *alternating* Prisoner's Dilemma game can affect the interaction to a considerable extent. For instance, if two Tit For Tat players engage in a Prisoner's Dilemma of the usual *simultaneous* kind, and if one of them defects by mistake, both players will subsequently cooperate and defect in turns. On the other hand, if two Tit For Tat players engage in an *alternating* Prisoner's Dilemma, and a unilateral defection occurs inadvertently, then the outcome will be an unbroken sequence of mutual defections.

In the alternating Prisoner's Dilemma, the roles of the two players in each round are asymmetrical. One of the player is the 'leader' (to use a game-theoretic expression) and able to decide what the outcome is going to be. In a single round, the leader obtains a higher payoff if he defects than if he cooperates. But if the leader cooperates, the other player receives a payoff which is higher than what he would receive if the leader defects; and this increment of his payoff is higher than the loss to the leader. If we consider one 'unit' of two consecutive rounds, and assume that each player finds himself during one of the two rounds in the role of leader, we find that their payoff, if both help each other, is larger than their payoff if both defect, but that if one player helps (when a leader) and the other does not (when it is his turn), then (a) the helper gets the lowest payoff—he has been had for a sucker—and (b) the defector obtains a payoff even higher than the *reward* for mutual cooperation, whereas (c) the average payoff for both players is lower than this reward. These are exactly the conditions defining the Prisoner's Dilemma. However, the interactions for the iterated game are now quite different.

In the alternating game, Pavlov is no longer error-correcting. A mistake by one of the players results in a run where each player defects every third round. Accordingly, evolutionary chronicles of the alternating game do not lead to the establishment of Pavlov. Much more successful is a strategy which always cooperates, except if it has been suckerpunched in the previous round. This strategy, which one may call *Firm But Fair*, is more tolerant than Tit For Tat, since it does forgive a defection by the opponent if this was in answer to the own defection. It is error-correcting.

In an interesting series of experiments, Claus Wedekind and Manfred Milinski have tested humans (more specifically, first year students) in both the simultaneous and the alternating Prisoner's Dilemma (see Wedekind and Milinski, 1996). In both cases the sub-

jects became more and more cooperative and successful as the experiment went on. It turned out that their strategies eventually clustered around two types which could be reasonably interpreted as versions of Pavlov respectively Firm But Fair. The players, however, were rather inflexible: they did not distinguish whether they were playing a simultaneous or an alternating game. This missing flexibility, according to Wedekind and Milinski, 'suggests either that people might have preferred niches, which contain either the simultaneous or the alternating situation, or that the game situation [in the experimental setup] did not offer the natural clues by which humans would recognise the mode and trigger their response conditionally... Pavlovian player appeared to gain more in the simultaneous game and [Firm But Fair]-players in the alternating game.' Apparently our propensity for cooperation has not been fine-tuned to a large extent. It should be noted that Milinski and Wedekind concentrated on finding strategies which depended only on what happened in the last round, and not on the players' standing. They did not consider Contrite Tit For Tat, for instance. But computer simulations show that Contrite Tit For Tat does very well (and even better than Firm But Fair) in the alternating Prisoner's Dilemma.

It seems rather difficult to test whether human subjects play Contrite Tit For Tat, by the way. The experimenter has to cause the subject to mis-implement a cooperative move in such a way that he feels guilt (rather than annoyance at the experimenters' interference).

# 6. THE SNOWDRIFT GAME

Another detail likely to plague ethologists on the look-out for mutual aid is the similarity between the Prisoner's Dilemma game and a close relative which has been termed the *Snowdrift* game by Robert Sugden: 'Suppose you are driving your car on a lonely road in winter, and you get stuck in a snowdrift, along with one other car. You and the other driver have both sensibly brought shovels. It is clear, then, that you should both start digging. Or is it? The other driver cannot dig his own way out of the drift without digging your way out too. If you think he is capable of doing the work on his own, why bother to help him?' (Sugden, 1986).

If we assume that it is better to do half the digging and get out of the snowdrift than to remain stuck, but that it is so important to get out of the snowdrift that each player would rather do all the digging himself than remain stuck, we obtain a rank ordering of the payoff given by $T>R>S>P$. This looks like the Prisoner's Dilemma, except that the payoffs $P$ and $S$ have been interchanged: the outcome for mutual defection is less good than that for unilateral cooperation. Better do the job by yourself than not get it done at all. This may seem, at first glance, a minor variation, but it is an entirely different ball-game, in fact.

As Sugden has pointed out, the Snowdrift game has the same structure as the much studied *Chicken* game, which played an essential role in providing an individual (rather than group) selection argument for the evolution of constraints in innerspecific fights (see Maynard Smith, 1982). In the Chicken game, the players have to decide whether to escalate a potentially dangerous conflict or whether merely to display (and retreat if the other player starts getting mean). If we label escalating by **D**, and displaying by **C**, we obtain the same rank ordering of the outcomes as with the Snowdrift game, in spite of the fact that the options in the two games have a totally different meaning: refusing help is certainly something else than offering an all-out fight.

Since the Chicken game has been thoroughly studied, we can use the corresponding results. It follows that for one round of the Snowdrift game, it is not necessarily the best to play **D**, i.e. to defect. The evolutionarily stable outcome will eventually depend on some asymmetry between the players, a seemingly arbitrary labelling of the players resulting from the more or less contingent history of the evolving population, reminiscent in human terms of a 'norm' or convention like driving on the left hand side on British roads. An asymmetry (in age, or strength, or ownership, etc.) decides which one of the two players opts for **C** and which for **D**. This outcome is entirely different from the seemingly similar Prisoner's Dilemma: we can no longer expect mutual cooperation, or a strategy like Pavlov or Contrite Tit For Tat.

This may shed some light on some recent controversy. A group of behavioural ecologists discovered that among lions, the principle of reciprocity seemed not to hold. 'That', it was claimed, 'throws a monkey wrench into the classic explanation for the evolution of cooperation in a selfish, dog-eat-dog (or lion-eat-gazelle) world' (see Morell, 1995).

Female lions typically live in prides of two to seven members, around which cluster their dependent offspring and a coalition of immigrant males. The males defend the pride against incursion by other males, whereas the females defend the territory against incursions by other females. Territorial incursions can be simulated by the playback of recorded roars, and these routinely elicit cooperative defense. In response, female lions pair off to approach what must seem to them an unknown and potentially dangerous enemy (see Heinson and Packer, 1995).

This is similar to the *predator inspection* whereby pairs of stickleback gingerly approach a pike to test its current mood and motivation. As long as the stickleback approach together, their risk is more than halved, and we can speak of cooperation. If one stickleback consistently lags behind and gains information by waiting in the wings, this is a clear case of defection. A beautiful series of experiments by Milinski has shown that this is an instance of the repeated Prisoner's Dilemma game (of the simultaneous sort, by the way), and that the stickleback use a strategy based on reciprocation (see Milinski, 1987).

The lions, in an apparently very similar situation, seem not to use the principle of reciprocity. Indeed, some of the lions turned out to be consistently laggards, others consistently bold defenders of the home range. The bold lions did not change their behaviour, even after having experienced in many 'rounds' the cowardly behaviour of the laggards. (Some of these laggards, incidentally, seem to use a conditional strategy, rushing forward only at the last minute.) More surprising than the diversity of behaviour among lions is its relative fixity: each lioness has her character—that of an intrepid leader attacking unconditionally, or that of a feet-dragging laggard shying away from every risk of an all too direct confrontation with a territorial intruder. The fact that the laggards are not punished, and that the leaders do not stop their advance, but just glance back, is a strong argument against a strategy based on Tit For Tat or Pavlov.

Does this imply the 'inadequacy of current theory to explain cooperation'? We believe that this is not the case. Indeed, remember that in the Prisoner's Dilemma, one of the essential prerequisites is that a player matched against a defector obtains a higher payoff if he defects than if he cooperates. Can we assume that this condition is met here? What happens to a lioness paired with a defector? If she cooperates by advancing boldly, she certainly runs a risk to be wounded in a fight against the intruder. But if the lioness retreats, this means that the pride is giving up its territory. As the authors of the lion study write, however, 'territory is essential for successful breeding'. It could well be that this is a graver threat to the reproductive success than the risk of squaring off against an intruder.

This is quite different from the scenario with the two stickleback approaching their pike. If a stickleback approaches all by himself, he risks a lot and gains relatively little; it is certainly better, in this case, to give up at an early stage of the approach, rather than learn about the motivational state of the pike at one's own expense. Just swimming quietly away does not mean giving up the territory. In contrast, if both lions defect, their situation becomes critical: 'Without a territory, a female lion has little chance of raising her cubs and so loses the chance to pass her genes on to the next generation—the bottom line of evolutionary success.'(see Heinson and Packer, 1995).

As soon as we admit this, we are faced with the Snowdrift game. Two lions belonging to the same pride and called upon defending their territory will know each other well and be aware of many asymmetries in their strength, age, status etc. It seems almost inescapable that the two lions find themselves in different roles, and hence will adopt different strategies for a Snowdrift game. Thus we ought not to be surprised when one of the lions is chickening out.

We have seen that seemingly slight variations in the game—leading to an alternating Prisoner's Dilemma, for instance, or to a Snowdrift game—can have remarkably different outcomes. There are other factors that are just as likely to affect the outcome in a fundamental way: for instance, if we assume that players can watch how their rivals do against each other, or if we allow that a player can break off a repeated game and start anew with another partner. These variants are likely to play an important role, especially in sophisticated communities—for instance, human tribes. The analysis of the repeated Prisoner's Dilemma game by means of theoretical or experimental mathematics provides at best a first approximation to the issue of mutual aid among humans, something comparable to the role of motion without friction in classical mechanics, which is valid for planetary motion, but hardly for anything more down-to-earth. Nevertheless, it provides an essential jump-off point for experimental research.

# 7. REFERENCES

Axelrod. R. (1984) *The Evolution of Cooperation*. Basic Books, New York (reprinted 1989 in Penguin, Harmondsworth).

Boerlijst, M., Nowak, M. and Sigmund, K. (1996), The logic of contrition, to appear in *Journal of Theoretical Biology*.

Boyd, R. (1989), Mistakes Allow Evolutionary Stability in the Repeated Prisoner's Dilemma Game. *Journal of Theoretical Biology* **136**, 47–56.

Frean, M.R. (1994), The Prisoner's Dilemma without synchrony, *Proceedings of the Royal Society London B*. **257**, 75–79.

Heinson, R. and Packer, C. (1995), Complex cooperative strategies in group-territorial african lions, *Science* **269** 1260–2.

Hofbauer, J. and Sigmund, K. (1988), *The Theory of Evolution and Dynamical Systems*, Cambridge UP.

Lindgren, K. (1991) Evolutionary phenomena in simple dynamics, in *Artificial life II* (ed. C.G. Langton et al), Santa Fe Institute for Studies in the Sciences of Complexity, X, 295–312.

Lindgren, K. (1996), Evolutionary Dynamics in Game-Theory Models, to appear in *The Economy as an Evolving, Complex System II*, ed. W. Brian Arthur et al.

Lindgren, K. and Nordahl, M.G. (1994), Evolutionary dynamics of spatial games, *Physica D* **75**, 292–309.

Maynard Smith, J. (1982) *Evolution and the theory of games*, Cambridge UP.

Milinski, M. (1987), Tit for tat in sticklebacks and the evolution of cooperation, *Nature* **325**, 434–435.

Morell, V. (1995), Cowardly lions confound cooperation theory, *Science* **269**, 1216–7.

Nowak, M. and Sigmund, K. (1992), Tit for tat in heterogeneous populations, *Nature* **355**, 250–2.

Nowak, M. and Sigmund, K. (1993), Win-stay, lose-shift outperforms tit-for-tat, *Nature*. **364**, 56–8.

Nowak, M. and Sigmund, K. (1994), The alternating Prisoner's Dilemma, *Journal of Theoretical Biology* **168**, 219–26.

Nowak, M., Sigmund, K. and El-Sedy, E. (1995), Automata, repeated games, and noise, *Journal of Mathematical Biology* **33**, 703–32.

Selten, R. (1975), Re-examination of the perfectness concept for equilibrium points in extensive games, *International Journal of Game Theory* **4**, 25–55.

Sinervo, B. and Lively, C.M. (1996), The rock-paper-scissors game and the evolution of alternative male strategies, *Nature* **380**, 240–243.

R. Sugden (1986), *The economics of Rights, Co-operation and Welfare*. Blackwell, New York.

Trivers, R. (1985), *Social Evolution*. Menlo Park CA, Benjamin Cummings.

Wedekind, C. and Milinski, M. (1996), Human cooperation in the simultaneous and the alternating Prisoner's Dilemma: Pavlov versus Generous Tit For Tat. *Proceedings of the National Academy of Science of the USA*, **93**, 2686–2689.

Weibull (1995) *Evolutionary Game Theory*, MIT Press, Cambridge, Mass.

Wu, J. and R. Axelrod (1995), How to cope with noise in the iterated Prisoner's Dilemma, *Journal of Conflict Resolution* **39**, 183–9.