

Albert, Max; Rusch, Hannes

**Working Paper**

## Indirect reciprocity, golden opportunities for defection, and inclusive reputation

MAGKS Joint Discussion Paper Series in Economics, No. 29-2013

**Provided in Cooperation with:**

Faculty of Business Administration and Economics, University of Marburg

*Suggested Citation:* Albert, Max; Rusch, Hannes (2013) : Indirect reciprocity, golden opportunities for defection, and inclusive reputation, MAGKS Joint Discussion Paper Series in Economics, No. 29-2013, Philipps-University Marburg, Faculty of Business Administration and Economics, Marburg

This Version is available at:

<https://hdl.handle.net/10419/93508>

**Standard-Nutzungsbedingungen:**

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

**Terms of use:**

*Documents in EconStor may be saved and copied for your personal and scholarly purposes.*

*You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.*

*If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.*

**MAGKS**



**Joint Discussion Paper  
Series in Economics**

by the Universities of  
**Aachen · Gießen · Göttingen  
Kassel · Marburg · Siegen**

ISSN 1867-3678

**No. 29-2013**

**Max Albert and Hannes Rusch**

**Indirect Reciprocity, Golden Opportunities for Defection,  
and Inclusive Reputation**

This paper can be downloaded from  
[http://www.uni-marburg.de/fb02/makro/forschung/magkspapers/index\\_html%28magks%29](http://www.uni-marburg.de/fb02/makro/forschung/magkspapers/index_html%28magks%29)

Coordination: Bernd Hayo • Philipps-University Marburg  
Faculty of Business Administration and Economics • Universitätsstraße 24, D-35032 Marburg  
Tel: +49-6421-2823091, Fax: +49-6421-2823088, e-mail: [hayo@wiwi.uni-marburg.de](mailto:hayo@wiwi.uni-marburg.de)

# Indirect Reciprocity, Golden Opportunities for Defection, and Inclusive Reputation

**Authors:** Max Albert<sup>1</sup>, Hannes Rusch<sup>1,2,†</sup>

<sup>1</sup> Chair of Behavioral and Institutional Economics, JLU Giessen, Licher Strasse 66, D-35394 Giessen, Germany

<sup>2</sup> Peter Löscher Chair of Business Ethics, TU München, Arcisstrasse 21, D-80333 Munich, Germany

<sup>†</sup> to whom correspondence should be addressed: hannes.rusch@uni-giessen.de, phone: +49-641-99-22-201 (front office), fax: +49-641-99-22-209

**Abstract:** In evolutionary models of indirect reciprocity, reputation mechanisms can stabilize cooperation even in severe cooperation problems like the prisoner's dilemma. Under certain circumstances, conditionally cooperative strategies ("cooperate iff your partner has a good reputation") cannot be invaded by any other strategy that conditions behavior only on own and partner reputation. Still, the evolutionary version of backward induction can lead to a breakdown of this kind of indirect reciprocity. Backward induction, however, requires strategies that count and then cease to cooperate in the last, last but one, last but two, ... game they play. These strategies are unlikely to exist in natural settings.

We present two new findings. (1) Surprisingly, the same kind of breakdown is also possible without counting. Opportunists using rare golden opportunities for defection can invade conditional cooperators. This can create further golden opportunities, inviting the next wave of opportunists, and so on, until cooperation breaks down completely. (2) Cooperation can be stabilized against these opportunists by letting an individual's initial reputation be inherited from that individual's parent. This 'inclusive reputation' mechanism can cope with any observably opportunistic strategy. Offspring of opportunists who successfully exploited a conditional cooperator cannot repeat their parents' success because they inherit a bad reputation, which forewarns conditional cooperators in later generations.

**Keywords:** evolutionary game theory; repeated prisoner's dilemma; backward induction; conditional cooperation; opportunism;

**JEL codes:** C73, D03, D64

## 1 Introduction

The study of problems of cooperation has been fascinating and connecting scholars from disciplines as diverse as economics, biology, anthropology, mathematics, psychology and philosophy, to name but a few, for decades now. Countless theoretical and empirical studies have been dedicated to solving parts of the puzzle of (human) cooperation (Axelrod and Hamilton, 1981; Clutton-Brock, 2009; Nowak, 2012; West et al., 2011; West et al., 2007). Various mechanisms that promote the evolution of cooperative behavioral strategies even in scenarios where individuals are unrelated and where defection is a strictly dominant strategy (as in, e.g., prisoner’s dilemmas, PDs) have been proposed and analyzed. The most prominent of these mechanisms are probably: direct reciprocity (Axelrod and Hamilton, 1981), network reciprocity (Nowak, 2006), kin selection (Lehmann et al., 2007; West et al., 2011), and indirect reciprocity (Brandt et al., 2007; Nowak and Sigmund, 2005). All these approaches are based on a process of (cultural or biological) evolution in which strategies are inherited.

In this paper, we focus on indirect reciprocity. This mechanism assumes that observing an individual’s past behavior is informative about its future behavior. Information about past behavior makes a strategy of conditional cooperation possible. In the simplest case, the information is public and transmitted in a binary variable called ‘reputation’, which can only be ‘good’ or ‘bad’, and which is used by conditional cooperators to identify other conditional cooperators. Individuals start their life with a good reputation. A conditional cooperator cooperates if and only if the partner has a good reputation. Deviations from this strategy, e.g., defecting against a partner with good reputation, reveal that an individual is not a conditional cooperator and, therefore, lead to a bad reputation. If there are sufficiently many interactions between individuals during one generation, conditional cooperators manage to cooperate mostly among themselves by identifying those individuals who use other strategies and stopping cooperation with them. This gives them an evolutionary advantage. If there are sufficiently many conditional cooperators to begin with, conditional cooperation spreads in the population until it is the predominant behavior, with defection occurring, if at all, only due to random events: mutations (i.e., random deviations of the offspring’s strategy from the parent’s strategy), trembles (i.e., random deviations of the action recommended by an individual’s strategy), or reporting errors (i.e., random errors in assigning, transmitting or identifying reputations).

However, indirect reciprocity leads to cooperation only if specific strategies are missing in the population. Missing strategies generate a ‘safe haven’ for conditional cooperation but also mean that ‘golden opportunities’ for defection are left unexploited. We borrow the term ‘golden opportunity’ from Frank (1988) but use it more broadly. In our sense of the term, a golden opportunity for defection is a situation where an individual’s expected gains from unilateral defection outweigh the expected costs even in an environment where conditional cooperation prevails. For instance, the last interaction in an individual’s life is always a golden opportunity for that individual. However, golden opportunities for defection could occur at any stage of life. Just assume that there is always a small probability that the current interaction differs from the usual ones, either because the probability that defection is observed is unusually low (the case considered by Frank, 1988), or because the gains from unilateral defection are unusually high (a case we consider in more detail below, see 3.4). In these cases, opportunists, who use the golden opportunity when, or if, it arises but otherwise mimic conditional cooperators, can invade a population of true conditional cooperators.

Opportunism raises two problems, a problem of explanation and a problem for the stability of cooperation. First, humans often seem to resist the temptation to exploit golden opportunities; the question, then, is how to explain this observation. This is Frank’s problem (Frank, 1988). Second, the success of opportunists may create new golden opportunities and, eventually, lead to a complete breakdown of cooperation. The most prominent case is (the evolutionary version of) backward induction: For simplicity, let individuals interact  $n$  times with randomly changing partners and have exact information about the number of the current interaction. A population of conditional cooperators can then be invaded by a mutant that is conditionally cooperative in the first  $n-1$  interactions but

defects in the last interaction, which is always a golden opportunity. Once this mutant has taken over, a new golden opportunity arises: defecting in the last but one interaction, i.e., interaction  $n-2$ , and so on, until cooperation breaks down completely. Counting interactions and defecting in one’s last interaction may not be a feasible strategy in real life. However, there are other, more realistic opportunistic strategies that also can invade a population of conditional cooperators and induce a breakdown of cooperation (see section 2 below).

Similar arguments may be raised in connection with at least some of the other mechanisms that have been proposed for explaining the evolution of cooperation. Missing strategies are obviously relevant for the success of ‘tit for tat’ in Axelrod’s tournaments (Axelrod, 1984), where participants did send in strategies for a finitely repeated PD with a large number of rounds. Since the number of rounds was unknown, it was virtually certain that the strategy of playing ‘tit for tat’ until the end and then defecting in the last round would be missing—not to speak of the many other tit-for-tat variants starting defection in the last but one round, last but two round, and so on, whose presence would be necessary for a complete breakdown of cooperation through backward induction.

It is, of course, almost inevitable that many strategies are missing in a given population. After all, the number of strategies in a finitely repeated PD explodes when the number of rounds goes up. For this reason, mechanisms relying on missing strategies are not necessarily implausible as an explanation for the stability of cooperation. Nevertheless, if players are intelligent, the strategies they employ may involve learning and may, therefore, discover golden opportunities in many environments. In such a scenario, many strategies would still be missing, but those that are present might be especially good at exploiting golden opportunities because they have been selected for recognizing them. Indeed, according to the “social intelligence hypothesis” (Whiten and van Schaik, 2007), the evolution of intelligence is largely driven by selection pressure in favor of more complicated strategies in social interaction, many of them “Machiavellian” like, e.g., tactical deception (McNally and Jackson, 2013) or, as discussed in this paper, exploitation of golden opportunities.

Conditional cooperators may find ways to identify opportunists before they strike. Frank argues that opportunists’ involuntary behaviors (e.g., facial expressions) might give them away (Frank, 1988). In our terminology, this would mean that strategies involving mimicry are missing. This implies that there would be selective pressure favoring better mimicry (exemplified, e.g., by the existence of good actors), which constantly threatens conditional cooperation.

In this paper, we consider an extended variant of indirect reciprocity that does not rely on missing strategies and, therefore, can cope with the problem of golden opportunities. It is based on the following central idea: Since opportunism is passed on, the behavior of an individual’s (biological or cultural) parent may be more informative than the individual’s own behavior. Using information about parent behavior means using information provided by the evolutionary process itself. As we show in this paper, conditionally cooperative strategies using this information can cope with the problem posed by golden opportunities for defection—even if no strategies are missing in the conditional cooperators’ environment.

## **2 Indirect reciprocity, missing strategies, and golden opportunities**

### *2.1 Modeling indirect reciprocity*

Models of indirect reciprocity are central to understanding how systems of reputation—in humans and, possibly, also in other species (Bshary and Grutter, 2006)—might have evolved. In these models, like in many others concerned with the evolution of cooperation, populations of unrelated, randomly matched individuals play one-shot PDs or asymmetrical ‘Giving Games’. These simple games are used because they possess the features of the ‘most stringent cooperative dilemma’ (Nowak, 2012)—they might be called ‘stress tests’ for mechanisms enabling cooperative behavior. Although individual defection is strictly dominant in these games, the resulting equilibrium of universal defection is inefficient: every individual would be better off under universal cooperation. One way of enabling

more cooperative strategies to evolve in this scenario is allowing individuals to make their behavior dependent on information about past behavior of the individual they are matched with, i.e., their opponent’s reputation.

For easier reference, we subsequently assume a specific protocol for the evolutionary process. In each period  $T = 0, 1, \dots$ , there exists a population of  $N \gg 2$  individuals, called a generation. For theoretical considerations, we consider the limiting case of an infinite population ( $N = \infty$ ), where stochastic effects vanish in the aggregate and the protocol yields the discrete-time replicator dynamics. All simulations in this paper involve finite populations (with  $N$  stated explicitly).<sup>1</sup> All individuals of a generation are randomly matched to play a sequence of  $n > 0$  two-player games, with a different partner in each game. For analytical feasibility and to ensure comparability with previous studies, we use the most stringent symmetrical cooperation problem in its simplest form, i.e., the symmetrical prisoner’s dilemma. We refer to the whole sequence of games played by a generation as the super-game. Each single two-person game is called a stage game. The number of stage games an individual has played is called the individual’s age. Individuals are always matched with partners of the same age. All individuals of a generation die simultaneously after their last stage game. The lifetime payoff of an individual is the sum of payoffs over its  $n$  stage games. The generation payoff is the sum of lifetime payoffs of a generation. Each generation of offspring consists, again, of  $N$  individuals. Thus, there are  $N$  slots in the each generation. The probability that a specific slot is filled by an offspring of a specific individual of the previous generation is equal to this individual’s lifetime payoff divided by the respective generation payoff.

The simplest way of incorporating reputation in this model is to assign each individual a single binary value, representing either a ‘good’ or a ‘bad’ reputation (a modeling approach introduced by Nowak and Sigmund, 1998a, 1998b). In the standard model, all newborn individuals have a good reputation (an assumption we modify in section 3). This reputation may change after each stage game. Individual reputations encode all that is remembered about the history of play.

Ohtsuki and Iwasa (2004, 2006) consider a space  $\Sigma$  of sixteen (stage game) strategies, which we denote as a four-letter string, for instance, CDCD. We refer to these strategies as “four-letter strategies”. The four letters indicate the individual’s action, C (‘cooperate’) or D (‘defect’), for each of the four situations GG, GB, BG, BB (in this sequence), where the first letter of a string like GB denotes the individual’s own reputation (in this case: good) and the second letter denotes the opponent’s reputation (here: bad). For instance, an individual whose strategy is CDCD cooperates, irrespective of its own reputation, if and only if meeting a partner with good reputation. In Ohtsuki and Iwasa’s models, an individual’s super-strategy is playing the same four-letter strategy in all stage games. In this paper, we later also allow individuals to use more flexible super-strategies which will be represented by eight letter strings, e.g., CDCD-DDDD, describing an individual’s stage game strategies in a first and a second specified situation, e.g., the first and the second stage game of a generation.

In addition to the number of different strategies present, a population is characterized by a reputation dynamic (Ohtsuki and Iwasa, 2004, 2006). This is a function  $\{GG, GB, BG, BB\} \times \{C, D\} \rightarrow \{B, G\}$  assigning a reputation, B or G, to an individual depending on the combination of reputations of both, the individuals and its opponent, and the individual’s action. The space  $\Delta$  of possible reputation dynamics has  $256 (= 2^8)$  elements. The reputation dynamic describes the way information about players’ observed behavior is evaluated. For instance, a reputation dynamic

---

<sup>1</sup> The Java source code of all simulation models is available from the corresponding author upon request. The theory for finite populations is complicated since all conceivable developments have positive probability. We nevertheless use finite-population simulations in order to show that our results are still highly probable in a stochastic environment.

$d \in \Delta$  satisfying  $d(GG,D) = B$  implies that a reputable individual which defects against another reputable individual is assigned a bad reputation.

Ohtsuki and Iwasa examine all combinations  $(d,p) \in \Delta \times \Sigma$  of reputation dynamics and strategies (Ohtsuki and Iwasa, 2004, 2006). They assume that each individual plays several games with different opponents. They find eight combinations  $(d_i,p_j)$ , called the ‘leading eight’, that were successful according to the following two criteria: (1) The level of cooperation in a population characterized by  $(d_i,p_j)$  is high even when trembles and reporting errors are possible. (2) Given the reputation dynamic  $d_i$ , the strategy  $p_j$  is an evolutionary stable strategy (ESS). According to Ohtsuki and Iwasa, these eight combinations can thus be called ‘desirable social norms’ (Ohtsuki and Iwasa, 2006).

|       | $d_k(GG,C)$ | $d_k(GG,D)$ | $d_k(GB,C)$ | $d_k(GB,D)$ | $d_k(BG,C)$ | $d_k(BG,D)$ | $d_k(BB,C)$ | $d_k(BB,D)$ |
|-------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| $d_1$ | G           | B           | B           | G           | G           | B           | B           | B           |
| $d_2$ |             |             | B           |             |             |             | G           |             |
| $d_3$ |             |             | B           |             |             |             | B           |             |
| $d_4$ |             |             | B           |             |             |             | G           |             |
| $d_5$ |             |             | G           |             |             |             | B           |             |
| $d_6$ |             |             | G           |             |             |             | G           |             |
| $d_7$ |             |             | G           |             |             |             | B           |             |
| $d_8$ |             |             | G           |             |             |             | G           |             |
| $d_9$ |             |             |             |             |             |             |             | B           |

**Table 1:** The leading eight reputation dynamics (plus a further dynamic,  $d_9$ , which is unforgiving)

The leading eight can be described as follows. All reputation dynamics  $d_i$ ,  $i = 1, \dots, 8$  satisfy  $d_i(GG,C) = G$ ,  $d_i(GG,D) = B$ ,  $d_i(BG,D) = B$ ,  $d_i(GB,D) = G$ , and  $d_i(BG,C) = G$  (see Table 1). There are two strategies, CDCD and CDCC. The eight leading pairs are  $(d_i, \text{CDCD})$ ,  $i = 1, 2, 4, 5, 6, 8$  and  $(d_i, \text{CDCC})$ ,  $i = 3, 7$ .

Since the strategy component of the leading eight pairs is almost constant, we speak of the leading eight reputation dynamics. These reputation dynamics all support a strategy of conditional cooperation: cooperate with those, and only with those, who have a good reputation. The reputation dynamics ensure that the number of conditional cooperators among those with a good reputation increases so that, in the long run, conditional cooperators cooperate predominantly among themselves. Note that this makes conditional cooperation an evolutionary stable strategy (ESS) only if there are small but positive rates of trembles or reporting errors. Without such errors, a population of conditional cooperators has no disreputable individuals, implying that unconditional cooperators would never be recognized as such and could, therefore, invade a population of conditional cooperators and spread by random drift.

The ‘forgiving’ nature of the leading eight— $d_i(BG,C) = G$ , that is, cooperation with a reputable individual restores a good reputation—is, however, not necessary for the maintenance of cooperation in the population. In combination with the ‘apologizing’ nature of the complementary strategies—that is, both strategies try to restore a good reputation by cooperating in the case BG—, and given that there are trembling hands or reporting errors, this feature just raises the rate of cooperation within a population.

Consider the unforgiving/unapologetic pair  $(d_9, \text{CDDD})$ , with  $d_9$  as described in Table 1. This pair will also support cooperation in equilibrium, although the level of cooperation will not be as high

as in the case of the leading eight if there are trembles and reporting errors. Since our point here is not the comparison of cooperation rates of whole populations, but rather the invasion of single populations by opportunistic strategies, we will restrict subsequent considerations to the dynamic  $d_g$ , the completely unforgiving, and thus strictest, reputation dynamic. As we will show, even under this dynamic, in which a single defection against a partner of good reputation will always lead to the exclusion of the perpetrator from the community of conditional cooperators, cooperation is vulnerable to invasion by opportunists. If the reputation dynamic is forgiving, this either does not change our results or (in the model of section 2.2.3) even favors opportunists, who can exploit forgiveness by defecting in the case of golden opportunities and regain their good reputation by apologizing if this is worth the costs. We will return to this point in the discussion (section 4).

## 2.2 *The problem of golden opportunities*

The modeling approach considered up to this point has delivered quite a number of valuable insights into the fundamental logic of reputation formation and its consequences for the maintenance of cooperative behavior. One striking observation, though, cannot be convincingly explained by the models considered so far: Why do (at least some) humans still cooperate even in the face of golden opportunities for defection (Frank, 1988)?

We define a golden opportunity for defection as a stage game where an individual with a good reputation can raise its expected lifetime payoffs by unilateral defection even if all other individuals are conditional cooperators. In order for defection to be profitable, it may be necessary, though, that the individual deviates from conditional cooperation in subsequent stage games, for instance, by defecting once it has lost its good reputation.

The results obtained by previous studies of indirect reciprocity fall under our category of missing strategies: strategies detecting and grasping golden opportunities are assumed to be absent. If, however, these strategies are present, all nine reputation dynamics fail to maintain cooperation under many circumstances. We will demonstrate this for three scenarios here: (i)  $n = 1$ ; (ii)  $n = 2$  with counting; (iii)  $n = 2$  without counting, but with golden opportunities that appear with a given probability  $g$ .

### 2.2.1 *Reputation is—obviously—useless for $n = 1$*

The simplest case is the case where every individual plays only one stage game per lifetime ( $n = 1$ ). In this setup, the set  $\Sigma$  of the sixteen four-letter strategies is the complete strategy space. With only one game to play, each game provides a golden opportunity for defection: since the reputation acquired by the individual in this game is irrelevant, defection is costless. Trivially, none of the reputation dynamics considered so far can maintain cooperation here. This is the problem of golden opportunities in its simplest and sharpest form. Any mechanism maintaining cooperation in this case should also be able to solve all cooperation problems involving larger numbers of stage games.

### 2.2.2 *Reputation cannot prevent evolutionary backward-induction when counting is possible*

Let us now assume that each individual plays twice ( $n = 2$ ). Since each individual plays twice, the space of super-strategies is much larger than the  $\Sigma$  of the sixteen four-letter strategies. We assume, as before, that individuals do not remember histories of play. In addition, we exclude strategies that condition actions in the second stage game on the opponent’s reputation in the first stage game. The resulting super-strategy space is  $\Sigma \times \Sigma$ , where an individual’s behavior is described by an eight-letter strategy  $p-q$ , i.e., a pair of four-letter strategies, with  $p$  for the first and  $q$  for the second stage game. This strategy space is already larger than necessary; specifically, given that all individuals start with a good reputation, only the first letter of  $p$ , which specifies the action in case GG, is relevant.

Assume that all individuals play the unapologetic variant of conditional cooperation CDDD in both stage games. Consider an individual with a good reputation. The individual’s behavior in the first



stage game maintains the individual’s good reputation. In the second stage game, this individual is unable to exploit the golden opportunity for defection that arises if it meets a partner with good reputation: it will again cooperate while the optimal behavior in the second game would, of course, be defection under all circumstances, that is, DDDD. The more flexible super-strategy CDDD-DDDD would successfully invade a population playing CDDD-CDDD, thereby making defection in the first stage game a golden opportunity. If all super-strategies of the type  $p$ - $q$  were present, only DDDD-DDDD would prevail, and the cooperation rate would be zero (except, possibly, for cooperation caused by trembles). Adding more super-strategies cannot change this result. This reasoning generalizes to any finite number  $n > 2$  of stage games: if individuals can ‘count’ and all super-strategies are present, the evolutionary analog of backward induction ultimately results in the breakdown of conditional cooperation, i.e., indirect reciprocity.

### 2.2.3 Golden opportunities without counting, $n = 2$

It might be argued, however, that identifying the last stage game of the game of life is a difficult task for individuals under all realistic circumstances. But even then, golden opportunities for defection may exist and lead to a breakdown of cooperation. As an example, we consider the case  $n = 2$ .

Assume that individuals play the same stage game twice. Each of the two stage games begins with the toss of a biased coin that determines whether payoffs are high or low. Under both conditions, low payoffs and high payoffs, individuals play a symmetric PD. Low payoffs are  $T > R > P > S$  with  $R > (T + S)/2$ , where the last condition ensures that cooperation is efficient. High payoffs are equal to low payoffs times  $a$  with  $a > 1$ . The high-payoff condition emerges with probability  $g < 0.5$  independently for each stage game.

Since individuals cannot count, they use the same strategy in each stage game. However, the super-strategy distinguishes between payoff conditions. For simplicity, we write these super-strategies again as eight letter strings, e.g., CDDD-DDDD, where now the first four-letter strategy describes behavior with low payoffs and the second describes the behavior with high payoffs. Conditional cooperators, then, are defined by CDXY-CDXY, that is, conditional cooperation under both payoff conditions, where X and Y each stand for C or D.

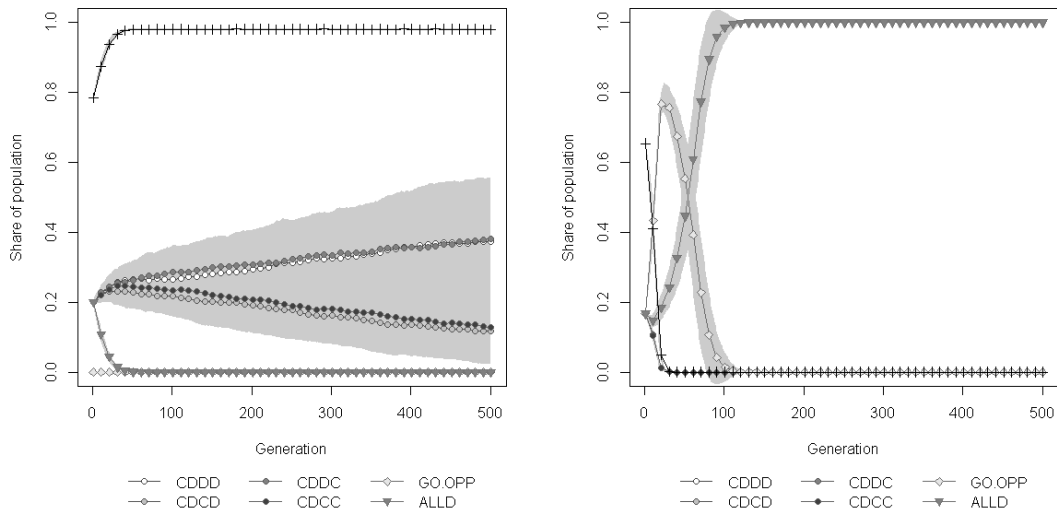
We put a further restriction on payoffs and assume that  $T-R < R-P$ , which implies that all-out defectors, described by DDDD-DDDD, cannot invade a population of conditional cooperators: relative to conditional cooperators, all-out defectors win  $T-R$  by defection in their first stage game but lose  $R-P$  in their second stage game.

Nevertheless, depending on parameter values, it is possible that the opportunistic strategy CDDD-DDDD, which conditionally cooperates under low payoffs but defects under high payoffs, invades a population of conditional cooperators, only to be displaced, once it has become established, by all-out defectors. This means that, although all-out defectors are unable to displace conditional cooperators on their own, they can take over once opportunists have paved the way for them.

Figure 1 shows simulation results with a finite population for a scenario ( $n = 2$  with  $T=15$ ,  $R=12$ ,  $P=6$ ,  $S=1$ ,  $g = 0.33$  and  $a = 4$ ) where all-out defectors, while unable to invade when opportunists are missing (Figure 1, left panel:  $N=5,000$ ; 4,000 conditional cooperators, 1,000 all-out defectors), take over after opportunists have displaced the initially dominant conditional cooperators (Figure 1, right panel:  $N=6,000$ , 4,000 conditional cooperators, 1,000 opportunists, 1,000 all-out defectors). Both simulations include trembles and reporting errors: In each stage game, an individual trembles (i.e., erroneously deviates from its strategy) with probability  $\varepsilon_{Act}$ . After each stage game, an individual is assigned the wrong reputation with probability  $\varepsilon_{Report}$ . Error probabilities are independent. We have set the probability of at least one tremble during an individual’s life (i.e., over two stage games) to approximately 0.01 (which requires  $\varepsilon_{Act} \approx 0.005$ ); the same probability holds for at least one reporting error in an individual’s life (thus,  $\varepsilon_{Report} \approx 0.005$ ). The simulation shows a case where the high-payoffs condition provides a golden opportunity for defection, and where the presence of

opportunists who use the golden opportunity destabilizes cooperation just as in the case of backward induction, although counting has been ruled out.

In support of the simulations, we analyze the limiting case of  $a = 1$  without trembles and reporting errors for an infinite population in appendix A. Even though the golden opportunity in the case  $a = 1$  offers exactly the same payoffs as the ordinary stage game, opportunists using this opportunity can, under certain conditions, destabilize cooperation just as in the simulation. This result shows that even an opportunistic strategy that just uses a rare event as an excuse for defection may destabilize cooperation.



**Figure 1:** The left panel shows simulations where opportunists (labeled ‘GO.OPP’) are absent and conditional cooperators (‘CD\*\*’) can therefore resist the invasion of all-out defectors (‘ALLD’). The right panel shows simulations where opportunists are present and are able to drive out conditional cooperators, only to be replaced by all-out defectors themselves thereafter. Apart from the absence or presence of opportunists, all parameters were identical across simulations. The panels show the results of 100 simulation runs each; points indicate mean frequencies of the respective strategies; grey areas show two standard deviations around the mean. The crossed line indicates the overall percentage of conditional cooperators with good reputations after all stage games have been played in that generation.

A general theory of the conditions under which exploitation of golden opportunities destabilizes cooperation is beyond the scope of the present paper. Our point here is to show, by example, that cooperation based on indirect reciprocity is more sensitive to invasion by opportunists than one might conclude from the literature. Although all-out defectors cannot invade populations of conditional cooperators, quite simple strategies exploiting golden opportunities for defection might do so and reduce the level of cooperation sufficiently to pave the way for all-out defectors.

An obvious solution to the problem of golden opportunities would be a more flexible ‘morality’, i.e., reputation dynamic: if defection in a golden opportunity (GO) did not lead to a loss of reputation, the breakdown of cooperation in the GOs would not affect cooperation in the ordinary games. However, this does not seem to be the kind of morality we observe: high gains are usually not considered as a valid excuse for defection, at least if defection also imposes a high cost on one’s partner.

### 3 Inclusive reputation can prevent opportunists from taking over

#### 3.1 Previously unused evolutionary information

In the scenarios considered in section 2, each newborn individual starts with a clean record, that is, a good initial reputation. Parents’ bad reputations do not lead to a bad reputation of their offspring,

although every individual is very likely to behave as well or as badly as its parent. This means that the reputation dynamics considered so far ignored relevant information, information that is provided by the evolutionary process itself, i.e., through the connection of parent and offspring. As long as the reputations and actions of individuals are publicly known (which is assumed in all models of indirect reciprocity discussed here) and as long as the parent of an individual is known, there exists reliable information about the individual even before it acts for the first time. The probability that the offspring of a disreputable individual is not a conditional cooperator is very high if mutation rates, tremble rates, and reporting-error rates are not too high. This is especially relevant for the problem of opportunism because opportunists are characterized by the fact that their behavior is indistinguishable from conditional cooperation until a golden opportunity for defection arises. In their case, recourse to their parents’ behavior is the only chance for successful prediction.

It is, therefore, a very plausible assumption that a reputation dynamic uses this information, which means that acquired reputations are passed on from parent to offspring. We label this idea ‘inclusive reputation’ because the reliable tradition of reputations from parent to offspring adds a new perspective to the calculation of the expected payoffs of a specific strategy. Over evolutionary time, only those strategies will persist which maximize their inclusive fitness (Hamilton, 1964), and when reputation is passed on, this inclusive fitness crucially depends not only on the parent’s, but also on the offspring’s reputation. Unless there are good reasons to assume that an individual’s parent is typically unknown, inclusive reputation dynamics should be taken into consideration. Remarkably, these reputation dynamics solve the problem of defection in golden opportunities without relying on missing strategies.

### 3.2 *The fundamentals of inclusive reputation, $n = 1$*

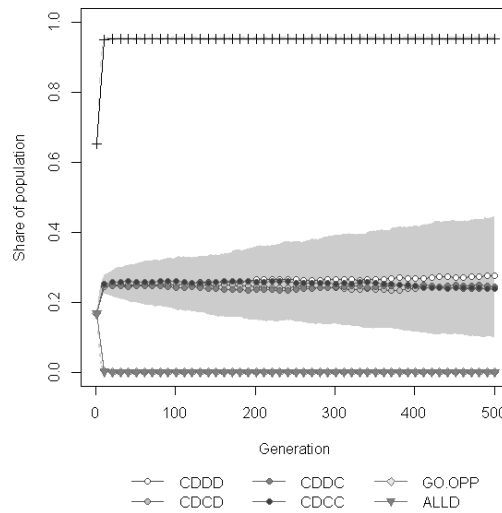
In this section, we focus again on the simplest cases of golden opportunities for defection, i.e.,  $n = 1$ . We consider only the unforgiving reputation dynamics  $d_0$  (see Table 1), which, in combination with inclusive reputation, supports a cooperative equilibrium. In this context, an unforgiving reputation dynamics implies that, unless a reporting error occurs, the offspring of a disreputable individual is always disreputable. Since we assume that each individual plays only one game, the sixteen functions from  $\{GG,GB,BG,BB\}$  to  $\{C,D\}$  describe the complete space of pure strategies. Let us assume, for simplicity, that there are neither mutations nor trembles nor reporting errors, and that the population is infinite. Then, with an inclusive reputation dynamic, the four conditionally cooperative strategies  $CDXY$  with  $X,Y \in \{C,D\}$  cannot be invaded by either all-out defectors or any kind of opportunists: Due to inclusive reputation, an individual that successfully exploits its good reputation by defecting against a conditional cooperator cannot pass on its advantage to its offspring: all of its offspring, through all generations, will be disreputable and will therefore not be able to exploit conditional cooperators. Moreover, in a population consisting mainly of conditional cooperators with good reputations, disreputable individuals will have less offspring than reputable conditional cooperators who almost always reap the benefits of mutual cooperation. For this simple reason, the inclusive fitness of those who defect against conditional cooperators is lower than that of the conditional cooperators themselves. Yet, from the perspective of an individual considering only its lifetime payoffs, each game is still a golden opportunity for defection.

The case  $n > 1$  with counting is even more favorable to conditional cooperators than the case  $n = 1$ . Since the last stage game is only one of several stage games, the excess return of last-game defectors relative to the payoff of conditional cooperators is smaller than in the case  $n = 1$ . Hence, inclusive reputation generally prevents evolutionary backward induction if mutants who count and defect in the last stage game appear.

Of course, without mutation rates, trembles and reporting errors, unconditional cooperators could invade and take over the population by random drift, as in the case without inclusive reputation. Once we add small rates of error, however, the four conditionally cooperative strategies are ESS again.

### 3.3 Golden opportunities without counting, with inclusive reputation, $n = 2$

The results of the case  $n = 1$  carry over to the model of section 2.2.3, where there are two stage games ( $n = 2$ ), individuals cannot count, and at each stage there is a probability that a golden opportunity presents itself. In section 2.2.3, we have shown that cooperation may break down completely because opportunists pave the way for all-out defectors. Again, inclusive reputation can stabilize conditional cooperation. Figure 2 shows the results of 100 simulation runs with exactly the same configuration as in Figure 1 (see 2.2.3), only that now reputation is inclusive.



**Figure 2:** Figure shows the results of 100 simulation runs for the same scenario as in Figure 1, right panel (i.e., a scenario where opportunists are present), only now with inclusive reputation enabled. Points indicate mean frequencies of the respective strategies; grey areas show two standard deviations around the mean; crossed line indicates overall share of conditional cooperators with good reputation after the two stage games.

In an infinite population, conditional cooperation with inclusive reputation is an ESS even for quite large golden-opportunity payoffs. Consider a single opportunist mutant appearing in a population of conditional cooperators. Since the mutant’s parent was a conditional cooperator, the mutant has a good reputation. However, once it encounters a golden opportunity, the mutant and almost all of its offspring have a bad reputation, which is inherited by almost all further offspring. Only by the occurrence of reporting errors can these individuals regain a good reputation. With a bad reputation, however, opportunists are equivalent to all-out defectors which cannot invade a population of conditional cooperators, not even without inclusive reputation, unless their initial number exceeds some threshold. Thus, unless the golden-opportunity payoffs are so large that the mutant’s offspring is above the threshold, conditional cooperation is an ESS.

When comparing the levels of conditional cooperators with good reputation in a generation, though, it can also be observed in the simulations that the benefits of inclusive reputation come at a price: Under inclusive reputation, trembles and reporting errors burden the offspring of conditional cooperators with bad reputations and, therefore, decrease cooperation rates as compared to a situation where conditional cooperators prevail and reputations are non-inclusive. When error rates are not too high, though, this only leads to a small decrease of overall cooperation levels.

### 3.4 Golden opportunities without counting for $n > 2$ and the power of inclusive reputation

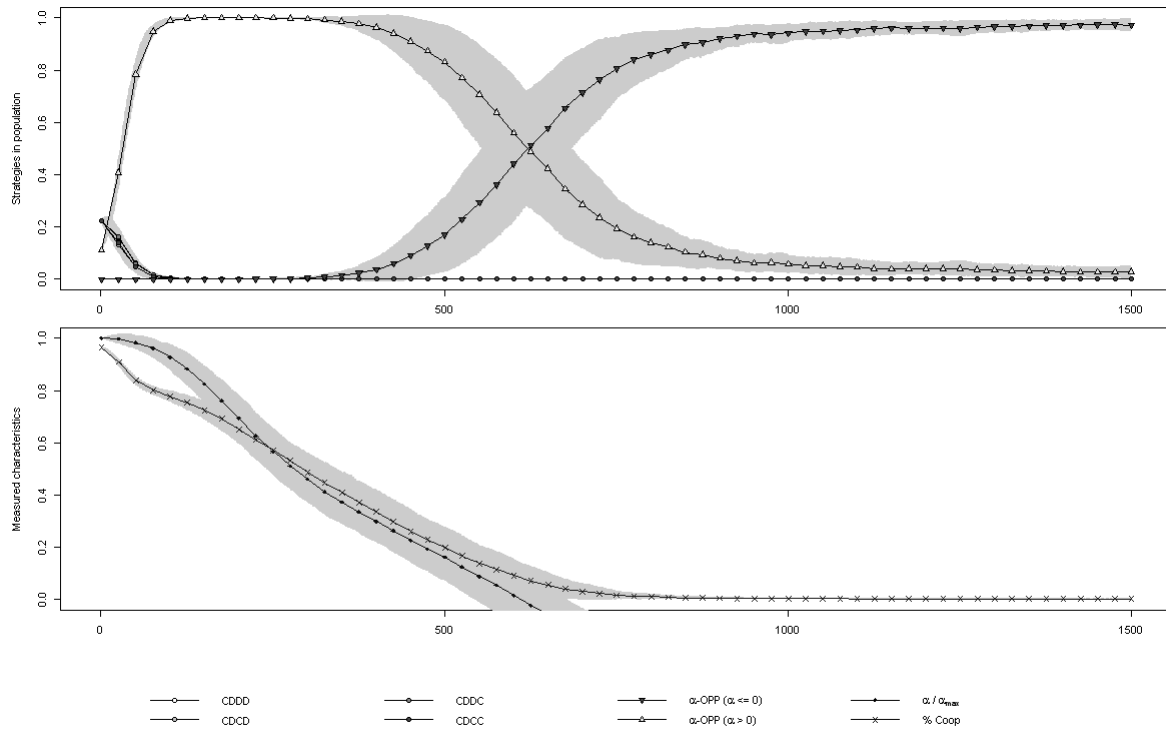
One might suspect that the effects we have demonstrated so far depend on the fact that we have only considered one or two stage games. Arguably, the standard, non-inclusive reputation mechanism

works best if there are many stages games so that conditional cooperators have time enough to identify defectors and opportunists and to enjoy the profits of cooperating among themselves. Of course, it is remarkable that inclusive reputation can solve problems where non-inclusive reputation fails, among them the problem of cooperation in the non-repeated PD. But, still, one might conjecture that the advantage of non-inclusive reputation might vanish when there are more stage games.

For this reason, we consider a second, more complex simulation, with the same basic protocol as before ( $N < \infty$  individuals; random matching for each stage game; ‘unforgiving’ reputation dynamic). However, we now consider  $n = 10$  stage games per generation. Reporting errors and trembles both appear with independent probabilities of ca. 0.001 per stage game, which implies a probability of approx. 0.01 each for a tremble or reporting error at least once in an individual’s life (which is the same life-time probability we assumed in 2.2.3).

Instead of just two payoff conditions, low and high, we now introduce a continuum of possible payoffs defined as follows. We use one specific symmetrical PD payoff matrix as a basis and multiply it by a random scalar  $(1 + a)$  with  $a \in [0, \infty)$  as the first move of each stage game. We compute  $a$  as  $a = (1 - u)^{-1/2} - 1$  where  $u$  is sampled from a uniform distribution with minimum 0 and maximum 1. This yields a Pareto distribution of  $a$  with minimum 0 and shape parameter 2. Thus, the payoffs of each stage game are given as  $(1+a) \cdot T$ ,  $(1+a) \cdot R$ ,  $(1+a) \cdot S$ , and  $(1+a) \cdot P$  analogously to 2.2.3. The distribution of  $a$  yields values in the proximity of 0 with high probability but infrequently also much higher values, thus creating stage games with very strong incentives to defect, i.e., golden opportunities, from time to time. We adapt individuals’ strategies to this new scenario by defining a threshold strategy called ‘alpha opportunism’. Before they make their move, individuals observe the payoffs or, equivalently, the specific realization of the scalar  $a$  that was chosen to modulate the payoffs. They then compare this value of  $a$  to their individual threshold value  $\alpha$ . If  $a \geq \alpha$  for a particular individual that individual defects; if  $a < \alpha$ , it behaves like a conditional cooperator. After a generation has completed playing the stage games, a new generation of individuals is compiled by copying  $N$  parental individuals’ strategies into the next generation with a probability proportional to their relative fitness, i.e., their overall payoff from all stage games divided by the sum of all individuals’ aggregate payoffs. During this process, small random mutations of  $\alpha$  are introduced by adding a small amount  $\varepsilon \cdot \varphi$  to  $\alpha$ , with  $\varphi$  normally distributed with mean 0 and variance 1 and  $\varepsilon = 0.125$ , each time when creating a new individual. If an individual’s  $\alpha$  value is 0 or less, that particular individual behaves like an all-out defector; however, because  $\alpha$  is still subject to mutation, the offspring of an individual with  $\alpha \leq 0$  may have a positive  $\alpha$  and, therefore, behave as a conditional cooperator for small enough payoffs.

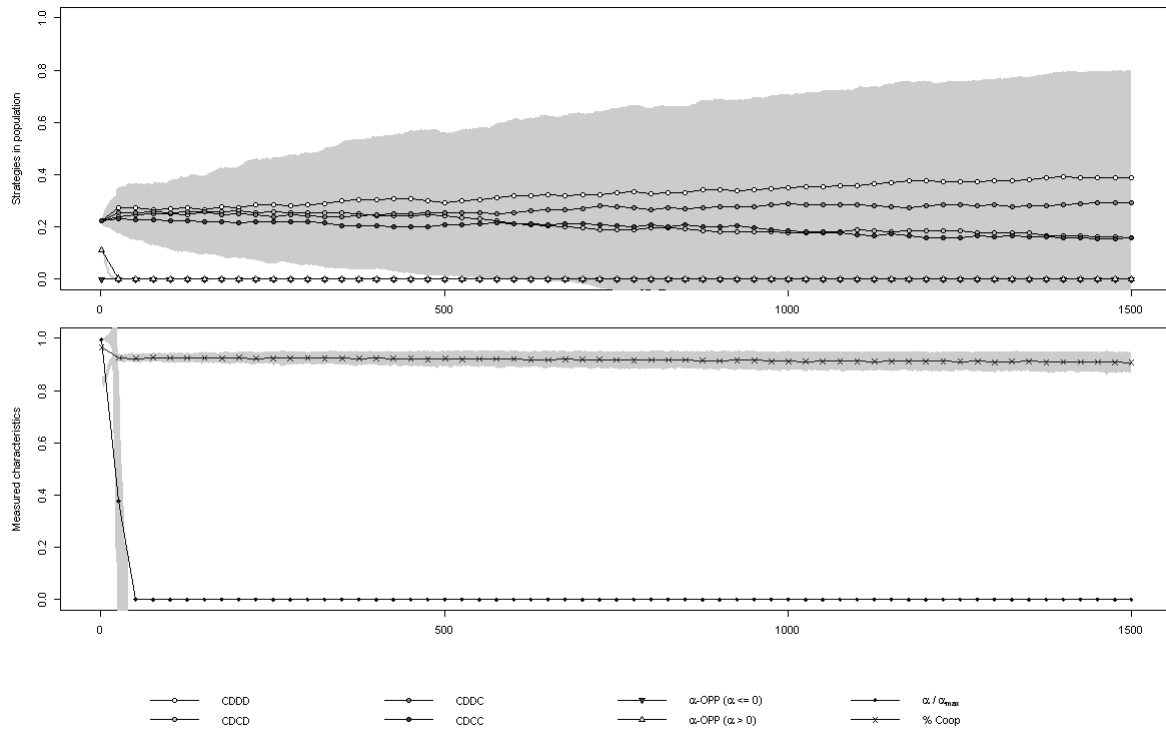
First, we consider the standard case without inclusive reputation. Experimenting with this simulation model shows that, under a wide range of circumstances, conditional cooperation is unable to prevent a breakdown of cooperation. For example, Figure 3 shows the evolutionary dynamics of 100 runs of the simulation for  $n = 10$ ,  $T=16$ ,  $R=9$ ,  $P=6$ ,  $S=1$  and  $N = 2,000$  (500 conditional cooperators of each variant, and 250 alpha opportunists with an initial  $\alpha$  of 6.0; error rates:  $\varepsilon_{Act} \approx \varepsilon_{Report} \approx 0.001$ , yielding independent life-time probabilities of approx. 0.01 for each kind of error).



**Figure 3:** An evolutionary arms race of alpha opportunists leads to a complete breakdown of cooperation. The upper panel shows the shares of conditional cooperators and opportunists in the population. Opportunists are counted in two distinct categories: opportunists with alpha levels greater than 0 (light triangles) and lower than 0 (dark triangles). The lower panel shows the mean value of all opportunists’ alphas, displayed as  $[\text{mean alpha}]/\text{max}(\text{alpha observed})$  for easier display, and the cooperation rate of the whole population. Points indicate mean values of the respective characteristic during 100 simulation runs, shaded areas indicate two standard deviations around the mean.

What happens is that, first, opportunists with high alpha thresholds drive out conditional cooperators by defecting when facing very high payoffs. Then, an evolutionary arms race leads to ever decreasing average levels of alpha, because among alpha opportunists those with alphas slightly lower than the average alpha level will have an advantage since, on average, they will defect a little earlier than the others and thus reap the payoff of one successful defection against a cooperator. After a period of falling average alpha levels, all-out defectors starts to appear again, in the form of opportunists with alphas equal to or below zero, as the average alpha level approaches zero, and take over the population almost completely. However, because the Pareto distribution of  $a$  implies only very small probabilities for values of  $a$  very close to 0, a small fraction of opportunists with positive alphas persists in the population, which, nonetheless, behave like all-out defectors almost all the time as well. This can be seen from the overall cooperation rate of the population, which converges to zero with falling alpha level. When switching from standard to inclusive reputation, however, conditional cooperators prevail right from the start and take over the population quickly, as shown in Figure 4.

We are well aware that the exact results are highly dependent on the specific values used for the initial composition of the population, initial alpha level, stage game payoffs, mutation frequency, error rates and number of stage games. However, the simulation exercise demonstrates that also for  $n > 2$  standard reputation dynamics cannot prevent an unraveling of cooperation under all circumstances and that the underlying problem can be fixed by introducing inclusive reputation, as long as error rates are not too high. The example shows that the intuitively quite evident notion of the power of inclusive reputation also stands more thorough tests, not only the case  $n = 2$ .



**Figure 4:** Same simulation scenario as in Figure 3, but with inclusive reputation enabled.

#### 4 Discussion

Although they foster cooperation, inclusive reputation dynamics are not necessarily ‘desirable social norms’. We grew up with the idea that persons should be judged on their own, not in the light of the behavior of their parents or grandparents. However, in the terms of Henrich et al. (2010), we are WEIRD (western, educated, industrialized, rich, democratic).

Moreover, our own moral preferences should not blind us to the fact that, until very recently, inherited reputation seems to have been very important in our own culture. It is not unreasonable to assume that the necessity of relying on inherited reputation is much reduced by modern institutions, which eliminate many golden opportunities for defection. As shown by Allen (2011), modern institutions became possible only through advances in technology like exact timekeeping, fast long-distance communication and reliable means of travel, which make it possible to determine whether failures in cooperative enterprises are due to bad luck or defection by a partner. Pre-modern institutions, in contrast to modern institutions, relied strongly on honor systems where the behavior of individuals determined the status of the family and, specifically, the individual’s offspring.

Indeed, as explanation for cooperative behavior, standard models of indirect reciprocity without inclusive reputation present us with two puzzles. Let us assume that these models sufficed for explaining cooperative behavior in humans. Why, then, do we find, in so many cultures, such a strong interest in the family background of other people? And why is it that keeping the wrong kind of company can ruin one’s reputation, even among WEIRD people? In both cases, it seems that trust is based on knowing where others come from. Inclusive reputation offers a potential explanation. The selection algorithm used in evolutionary modeling can either represent the genetic inheritance of behavioral traits between generations or the copying of successful behavioral strategies via social learning (by an ‘imitate the best’ rule), also between unrelated individuals. If behavioral strategies are, at least partially, genetically inherited, the desire to know the moral standing of a partner’s parents makes perfect sense. If, on the other hand, successful behavioral strategies are preferentially learned from others, the reputation of a person’s friends may also be a good predictor for a person’s behavior.

These two observations present difficulties for models of indirect reciprocity without inclusive reputation. Three other difficulties were already mentioned. First, the standard models assume that relevant and available information remains unused, namely, information about parents' reputations. This should be explained. Second, in the presence of more flexible strategies, the standard models predict that the maintenance of cooperation requires lower moral standards where using a golden opportunity to defect has no negative impact on a person's reputation. It seems to us, however, that this is actually not the case. And third, people are not always using golden opportunities for defection.

The results of this paper suggest that an environment that supports indirect reciprocity is very likely an environment that strongly favors inclusive reputation. Individuals who use information about opponents' parents' behavior employ a more complicated strategy and show, in this sense, a higher social intelligence. Thus, inclusive reputation belongs to the class of phenomena analyzed in models of 'Machiavellian intelligence', where selection pressure leads to more complex strategies in social interaction (Whiten and van Schaik, 2007).

In analogy to inclusive fitness, our approach adds a longer-term perspective to the discussion of indirect reciprocity. Those who use a golden opportunity for defection might gain quite tempting immediate benefits, but by revealing their true character they reduce their inclusive fitness—so much so that indirect reciprocity based on inclusive reputation can prevail even when golden opportunities for defection are present and no strategies are missing.

In this paper, we consider reputation only in relation to bilateral cooperation; moreover, we consider only an inclusive reputation dynamics that uses information about parent behavior. However, inclusive reputation may be more generally based on information about kin behavior. A kinship group may profit if one member signals the existence of a trait like cooperativeness or vengefulness in the whole group. Interesting trade-off problems between individual fitness and inclusive reputation might arise here. For example, the signal could be extremely costly to the individual while still being profitable to the group. Even if sending the signal kills the individual before it can produce any offspring, the benefit for its kin may be large enough for the trait to spread. This leads to all kinds of new problems; for instance, reputation becomes a public good among a kin group, which implies the possibility of free-riding strategies.

Inherited reputation may be important in real social interaction. However, empirical investigations in this direction are certainly needed, since the evidence we can currently refer to is mostly anecdotal (with the exception, possibly, of the historical evidence on honor systems mentioned above). Questions to investigate are: Does genealogical reputation information about persons or information about their preferred peers' reputations significantly influence other persons' cooperative attitudes towards them? Can traditions of good or bad reputations along chains of players who are allowed to advise their successors in cooperation games (a paradigm used, e.g., by Chaudhuri et al., 2006) be observed and manipulated experimentally? Empirical studies along these lines might allow us to establish empirically whether inherited reputations actually play a role for cooperative behavior.

We have, of course, focused on the strictest version of inclusive reputation, where the offspring of individuals violating the norms of indirect reciprocity can never restore the family reputation. One might inquire whether more lenient forms of inclusive reputation may be sufficient to foster cooperation. For instance, one could consider a reputation dynamic where an individual with an inherited bad reputation can acquire a good reputation by cooperating a few times with reputable individuals, thereby giving up what its parent gained by using a golden opportunity for defection. Thus, in analogy to the analyses by Ohtsuki and Iwasa (2004, 2006), such a forgiving reputation dynamics together with apologizing strategies might lead to higher levels of cooperation in the presence of randomly occurring action and reporting errors. Moreover, such more lenient reputation dynamics might be more in line with WEIRD moral preferences.



## Appendix A

We consider a scenario with two periods, a symmetrical PD with payoffs  $T > R > P > S$  in each period, and only two super-strategies, conditional cooperation (CDDD) and a super-strategy called opportunism. Opportunists decide, individually and before each stage game, whether they play CDDD or whether they switch from CDDD to all-out defection (DDDD). In each of these decisions, the probability of a switch is  $g \in (0,1)$ . Once an individual has switched to DDDD, it sticks to it. Of course, with probability  $(1-g)^2$ , an opportunist never switches.

All individuals start with a good reputation. We assume the unforgiving reputation dynamics  $d_0$ . With two periods, opportunists who switch in the first period defect against an individual with good reputation and therefore lose their own good reputation. If two opportunists with good reputations meet, both have not yet switched. The share of opportunists with a bad reputation in period  $j = 1, 2$  is denoted  $b_j$ , with  $b_1 = 0$  and  $b_2 = g$ .

If two opportunists with good reputations meet, switching could, in principle, be more or less correlated. In the case where switching results from observing high payoffs (see section 2.2.3), switching is perfectly correlated. We stick with this assumption. It implies that opportunists with good reputation playing against each other each have an expected payoff of  $(1-g)R + gP$ . With completely uncorrelated switching, the expected payoff in this situation would be  $g^2P + (1-g)^2R + g(1-g)(T+S)$ . Obviously, then, the expected payoff in the correlated case is higher than in the uncorrelated case if and only if  $P-S > T-R$ , that is, in the language of Rapoport (1967), if the “fear” factor  $P-S$  exceeds the “greed” factor  $T-R$ . This condition is stronger than the usual restriction on payoffs,  $R-P > T-R - (P-S)$  (“cooperators’ gain” exceeds “greed” minus “fear”) or, equivalently,  $R > (T+S)/2$ , which ensures that mutual cooperation pays more than taking turns in playing C against D. The condition  $P-S > T-R$  reappears below when we state sufficient conditions for the scenario behind Figure 1 in the text.

Let the population share of conditional cooperators be  $q$ . The expected payoff of a conditional cooperator in period  $j$ , then, is  $CC_j = qR + (1-q)b_jP + (1-q)(1-b_j)gS + (1-q)(1-b_j)(1-g)R$ , yielding a total expected payoff

$$(A1) \quad CC = CC_1 + CC_2 = [2q + (1-q)(1-g)(2-g)](R-P) + (1-q)g(2-g)(P-S).$$

The expected payoff of an opportunist with good reputation in the second period is  $OPP_2 = qgT + q(1-g)R + (1-q)b_2P + (1-q)(1-b_2)[gP + (1-g)R]$ . With a bad reputation, the payoff is always  $P$ . The total expected payoff, then, is

$$(A2) \quad OPP = qg(T+P) + q(1-g)(R+OPP_2) + 2(1-q)gP + (1-q)(1-g)(R+OPP_2) \\ = 2qg(T-R) + [q(2-g) + (1-q)(1+(1-g)^2)](R-P).$$

The difference in total expected payoffs between an opportunist and a conditional cooperator,  $\Delta = OPP - CC$ , is

$$(A3) \quad \Delta(q) = qg(2-g)(T-R) - g[q + (1-q)(1-g)^2](R-P) + (1-q)(2-g)g(P-S)$$

where  $\Delta'(q) = g(2-g)[T-R - g(R-P) - (P-S)]$ . Thus,  $\Delta$  decreases with  $q$  if and only if  $T-R < g(R-P) + (P-S)$ . Under the latter condition, opportunists can invade and completely take over a population of conditional cooperators if  $\Delta(1) > 0$  (because then their advantage increases as they spread). We find  $\Delta(1) > 0$  if and only if  $T-R > (R-P)/(2-g)$ . Thus, a sufficient condition for a successful takeover is

$$(A4) \quad (R-P)/(2-g) < T-R < g(R-P) + P-S.$$

This requires  $(R-P)/(2-g) < g(R-P) + P-S$  or, equivalently,  $(P-S)/(R-P) > (1-g)^2/(2-g)$ .

To summarize: (A4) is a sufficient condition for opportunists to crowd out conditional cooperators. This condition can be consistent with  $T-R < R-P$ , which prevents all-out defectors from invading a population of conditional cooperators.

We derive a further condition under which all-out defectors can invade and take over the population of opportunists. The following considerations, which are largely analogous to the preceding analysis, are simplified if we use a stronger sufficient condition instead of (A4):

$$(A5) \quad (R-P)/(2-g) < T-R < P-S.$$

Let  $s$  be the population share of the opportunists and  $1-s$  the population share of all-out defectors (assuming that conditional cooperators have died out completely). An opportunist’s expected payoff is

$$(A6) \quad OPP^* = s(1-g)[1+(1-g)^2](R-P) - (1-s)(1-g)(P-S),$$

while an all-out defector’s expected payoff is

$$(A7) \quad DD = s(1-g)(T-R) + s(1-g)(R-P)$$

This yields an expected payoff difference  $\Pi = OPP^* - DD$  depending on  $s$ :

$$(A8) \quad \Pi(s) = s(1-g)(T-R) - s(1-g)^3(R-P) + (1-s)(1-g)(P-S).$$

We find  $\Pi'(s) = (1-g)[T-R - (1-g)^2(R-P) - (P-S)]$ , which is negative under condition (A5). Therefore, all-out defectors have an advantage for all  $s \in [0,1]$  if and only if  $\Pi(1) > 0$  or, equivalently,

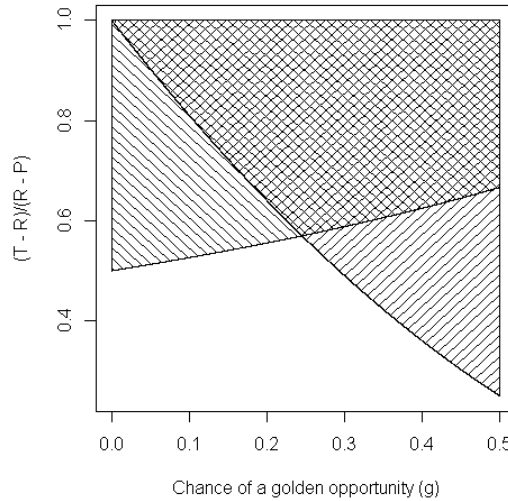
$$(A8) \quad (T-R)/(R-P) > (1-g)^2.$$

Combining (A5) and (A8) yields

$$(A9) \quad (T-R)/(R-P) > \max\{(1-g)^2, 1/(2-g)\} \text{ and } (T-R)/(P-S) < 1.$$

Let  $(T-R)/(R-P) < 1$ , implying that all-out defectors cannot, on their own, invade a population of conditional cooperators. Condition (A9), then, is a sufficient condition for the following scenario. Let the population be initially dominated by conditional cooperators, with small population shares for opportunists and all-out defectors. If (A9) holds, opportunists crowd out conditional cooperators but are, in their turn, crowded out by all-out defectors.

The condition  $(T-R)/(R-P) > \max\{(1-g)^2, 1/(2-g)\}$  is illustrated in figure 5.



**Figure 5:** The downward-sloping curve is  $(1-g)^2$ . The upward-sloping curve is  $1/(2-g)$ . Both curves intersect at  $g \approx 0.245$ . The shaded areas over the curves indicate the values of  $(T-R)/(R-P)$  satisfying conditions (1) and/or (2) for given values of  $g$ .

### Acknowledgements

For valuable suggestions and comments, we thank Matthias Greiff (Giessen), Bernd Lahno (Frankfurt), Georg Nöldeke (Basel), Jorge Peña (Basel), Matthias Uhl (Munich), and Eckart Volland (Giessen).

## References

- Allen, D.W., 2011. *The institutional revolution: Measurement and the economic emergence of the modern world*. Chicago: University of Chicago Press.
- Axelrod, R., Hamilton, W.D., 1981. The evolution of cooperation. *Science* 211, 1390-1396.
- Axelrod, R.M., 1984. *The evolution of cooperation*. New York: Basic Books.
- Brandt, H., Ohtsuki, H., Iwasa, Y., Sigmund, K., 2007. A survey of indirect reciprocity. In: Takeuchi, Y., Iwasa, Y., Sato, K. (Eds.). *Mathematics for ecology and environmental sciences*, Berlin: Springer, 21-49. (doi: 10.1007/978-3-540-34428-5\_3)
- Bshary, R., Grutter, A.S., 2006. Image scoring and cooperation in a cleaner fish mutualism. *Nature* 441, 975-978. (doi: 10.1038/nature04755)
- Chaudhuri, A., Graziano, S., Maitra, P., 2006. Social learning and norms in a public goods experiment with inter-generational advice. *Review of Economic Studies* 73, 357-380. (doi: 10.1111/j.1467-937X.2006.0379.x)
- Clutton-Brock, T.H., 2009. Cooperation between non-kin in animal societies. *Nature* 462, 51-57.
- Frank, R.H., 1988. *Passions within reason: The strategic role of the emotions*. New York: Norton.
- Hamilton, W.D., 1964. The genetical evolution of social behaviour I&II. *Journal of Theoretical Biology* 7, 1-52. (doi: 10.1016/0022-5193(64)90038-4)
- Henrich, J., Heine, S.J., Norenzayan, A., 2010. The weirdest people in the world? *Behavioral and Brain Sciences* 33, 61-83. (doi: 10.1017/S0140525X0999152X)
- Lehmann, L., Keller, L., West, S.A., Roze, D., 2007. Group selection and kin selection: Two concepts but one process. *Proceedings of the National Academy of Sciences* 104, 6736-6739. (doi: 10.1073/pnas.0700662104)
- McNally, L., Jackson, A.L., 2013. Cooperation creates selection for tactical deception. *Proceedings of the Royal Society B* 280, 20130699. (doi: 10.1098/rspb.2013.0699)
- Nowak, M.A., 2006. Five rules for the evolution of cooperation. *Science* 314, 1560-1563. (doi: 10.1126/science.1133755)
- Nowak, M.A., 2012. Evolving cooperation. *Journal of Theoretical Biology* 299, 1-8. (doi: 10.1016/j.jtbi.2012.01.014)
- Nowak, M.A., Sigmund, K., 1998a. Evolution of indirect reciprocity by image scoring. *Nature* 393, 573-577. (doi: 10.1038/31225)
- Nowak, M.A., Sigmund, K., 1998b. The dynamics of indirect reciprocity. *Journal of Theoretical Biology* 194, 561-574. (doi: 10.1006/jtbi.1998.0775)
- Nowak, M.A., Sigmund, K., 2005. Evolution of indirect reciprocity. *Nature* 437, 1291-1298. (doi: 10.1038/nature04131)
- Ohtsuki, H., Iwasa, Y., 2004. How should we define goodness? Reputation dynamics in indirect reciprocity. *Journal of Theoretical Biology* 231, 107-120. (doi: 10.1016/j.jtbi.2004.06.005)
- Ohtsuki, H., Iwasa, Y., 2006. The leading eight: Social norms that can maintain cooperation by indirect reciprocity. *Journal of Theoretical Biology* 239, 435-444. (doi: 10.1016/j.jtbi.2005.08.008)
- Rapoport, A., 1967. A note on the "index of cooperation" for Prisoner's Dilemma. *Journal of Conflict Resolution* 11, 100-103. (doi: 10.1177/002200276701100108)
- West, S.A., El Mouden, C., Gardner, A., 2011. Sixteen common misconceptions about the evolution of cooperation in humans. *Evolution and Human Behavior* 32, 231-262. (doi: 10.1016/j.evolhumbehav.2010.08.001)
- West, S.A., Griffin, A.S., Gardner, A., 2007. Social semantics: altruism, cooperation, mutualism, strong reciprocity and group selection. *Journal of Evolutionary Biology* 20, 415-432. (doi: 10.1111/j.1420-9101.2006.01258.x)
- Whiten, A., van Schaik, C.P., 2007. The evolution of animal 'cultures' and social intelligence. *Philosophical Transactions of the Royal Society B* 362, 603-620. (doi: 10.1098/rstb.2006.1998)