

This document is published at:

Antonioni, A., Sánchez, A. y Tomassini, M. (2016). Cooperation Survives and Cheating Pays in a Dynamic Network Structure with Unreliable Reputation. *Scientific Reports*, 6, 27160.

DOI: <https://doi.org/10.1038/srep27160>



© 2016 Nature Publishing Group. This article is licensed under a under a [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/). Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.

# SCIENTIFIC REPORTS



OPEN

## Cooperation Survives and Cheating Pays in a Dynamic Network Structure with Unreliable Reputation

Received: 04 March 2016

Accepted: 13 May 2016

Published: 02 June 2016

Alberto Antonioni<sup>1,2,3</sup>, Angel Sánchez<sup>2,3,4</sup> & Marco Tomassini<sup>1,2</sup>

In a networked society like ours, reputation is an indispensable tool to guide decisions about social or economic interactions with individuals otherwise unknown. Usually, information about prospective counterparts is incomplete, often being limited to an average success rate. Uncertainty on reputation is further increased by fraud, which is increasingly becoming a cause of concern. To address these issues, we have designed an experiment based on the Prisoner's Dilemma as a model for social interactions. Participants could spend money to have their observable cooperativeness increased. We find that the aggregate cooperation level is practically unchanged, i.e., global behavior does not seem to be affected by unreliable reputations. However, at the individual level we find two distinct types of behavior, one of reliable subjects and one of cheaters, where the latter artificially fake their reputation in almost every interaction. Cheaters end up being better off than honest individuals, who not only keep their true reputation but are also more cooperative. In practice, this results in honest subjects paying the costs of fraud as cheaters earn the same as in a truthful environment. These findings point to the importance of ensuring the truthfulness of reputation for a more equitable and fair society.

In present-day networked society a great many social and commercial interactions take place on internet<sup>1</sup>. In most instances, such interactions involve people who know each other only through an online identity<sup>2</sup>, without any connection whatsoever in the physical world. This is the case, for example, of internet platforms allowing private sales or exchanges among individuals<sup>3,4</sup>. In a different but related setting, a host of internet services and physical businesses (e.g., restaurants, hotels, etc.) rely on their online reputation to attract and keep their customers. Key to all these interactions is the reliability of the knowledge on the interaction counterpart, an issue that generates enormous concern these days due to the mounting evidence of fraud<sup>5</sup>. Consumer review websites such as Yelp or TripAdvisor use sophisticated analysis tools to remove (positive or negative) fake reviews; in fact, a whole new technical sub-field called *Online Reputation Management* dealing with how to detect, avoid, and eliminate fake reviews in online sites has recently arisen<sup>6,7</sup>. These concerns are even more pressing when personal identities, whose reliability is not externally checked, are the only available information about a possible interaction partner.

In this paper, we address this issue by framing the question in a simplified environment as a dyadic Prisoner's Dilemma (PD)<sup>8,9</sup> which lends itself to an experimental approach. Indeed, in online exchanges such as those described above, the best joint outcome obtains when both parties involved meet their end of the bargain, but both of them have clear incentives to cheat. In this situation, the game-theoretical prediction picks out defection as the rational choice in this game, but cooperation is often observed in our society and, in particular, in online exchanges. To explain this apparent paradox, several mechanisms have been proposed (see<sup>10</sup> for a recent review), most of which rely on some form of positive assortment between cooperators<sup>11</sup>, i.e., cooperators interact with individuals of similar behavior and avoid cheaters. In this context, both theoretical models<sup>12–16</sup> and recent experiments with human subjects<sup>17–21,23</sup> have established that cooperation may evolve to a remarkable degree when individuals control with whom they interact. Crucially for our present purposes, the process depends on

<sup>1</sup>Faculty of Business and Economics, University of Lausanne, 1015 Lausanne, Switzerland. <sup>2</sup>Grupo Interdisciplinar de Sistemas Complejos (GISC), Departamento de Matemáticas, Universidad Carlos III de Madrid, 28911 Leganés, Madrid, Spain. <sup>3</sup>Instituto de Biocomputación y Física de Sistemas Complejos (BIFI), Universidad de Zaragoza, 50018 Zaragoza, Spain. <sup>4</sup>Institute UC3M-B5 for Financial Big Data, Universidad Carlos III de Madrid, 28911 Leganés, Madrid, Spain. Correspondence and requests for materials should be addressed to A.A. (email: alberto.antonioni@unil.ch)

the availability of information on current and possible partners, which subjects then use to evaluate reputation<sup>20,21,23–29</sup> and to decide on their connections. It is then clear that the cooperation-promotion mechanism can act only if reputation scores are truthful, really reflecting the actual individual's record and are not manipulated in any way.

Here we contribute to the research on fake reputation and its effects by carrying out a controlled experiment<sup>30</sup> using a PD experiment with the possibility for participants to modify their behavior record by paying a cost. Such cost represents the effort that has to be done to pay or convince somebody else to alter our reputation in order to appear better than that we actually are or to decrease the reputation of a competitor. We could also have considered a cost-free alteration of one's reputation but, this fact being common knowledge in the experiment, it would have made the concept of a reputation almost useless. Our setup allows us to study whether having individuals with fake reputations around can undermine the evolution of cooperation and the success of dyadic online exchanges. This experimental approach, which to the best of our knowledge has not been attempted before, complements nicely the work carried out from the viewpoint of analyzing fraud evidence and associated behaviors in real systems<sup>31</sup>. As we will see, our results provide new insights on how people behave when they have the possibility to cheat and what are the consequences for the group: Thus, we will show that cooperation is not suppressed by the presence of individuals with fake reputation, but the society splits in two groups, one of them exploiting the other by cheating, leading to a sizeable increase in global inequality.

### Experimental setup

In our experimental sessions, seven groups of twenty subjects connected in a social network played a Prisoner's Dilemma (PD) game<sup>8,9</sup> with their neighbors. In this two-person game, players must decide whether to cooperate (C) or to defect (D) and, similarly to several recent experimental settings (e.g.<sup>17–19,21,23</sup>), the chosen action is the same with all neighbors. Note that if actions could be chosen independently for each neighbor the network disappears, and the system is simply a collection of independent pairwise games. If both agents cooperate, each receives a payoff  $R$ . If one defects and the other cooperates, the defector receives  $T$  and the cooperator receives the payoff  $S$ . If both defect, each receives  $P$ . Since  $T > R > P \geq S$ , defection is a dominant strategy and a rational payoff-maximizing player will choose to defect, although mutual cooperation yields a higher collective payoff, whence the dilemma. Subjects played a weak PD game ( $P = S$ ) with their immediate neighbors with  $T = 10$ ,  $R = 7$ ,  $P = 0$ , and  $S = 0$ . Payoff values are the same as those used in<sup>21,22</sup>, where it was shown that when the game is played on a static network cooperation decays, while the possibility to rewire links allows for its emergence and stability when information about past actions of others, which amounts to their reputation, is available. The initial set of connections between the participants was chosen to be a regular lattice of degree 4. Participants played 30 rounds of the sequence described below, although this exact number was unknown to them; they were only told that they would play for a number of rounds between 20 and 50 and without showing them the current round number.

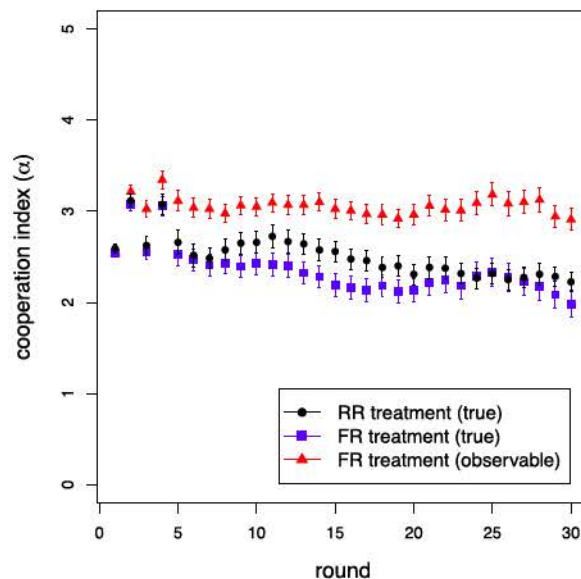
Here, the reputation of a player is expressed through a *cooperation index*  $\alpha$  which is the number of times the player has cooperated in the last five moves, thus  $\alpha \in [0, 5]$ . We considered two treatments: a baseline one, called *Real Reputation* (RR) in which the cooperation index cannot be manipulated, and a modified one in which participants were informed that all of them were allowed to vary their cooperation index by paying a cost, called *Fake Reputation* (FR). At the beginning, all players receive an initial  $\alpha$  of 3 based on the actions sequence CDCDC. Note that this form of reputation is related to but different from the one used in<sup>21,23</sup>. While in those earlier studies explicit past choices of each player were available to all others, in our experiment, there is some uncertainty about the current behavior of a player even in the RR treatment. This uncertainty comes about because only the number of cooperative actions of the current first neighbors and potential partners is known, but not their order. In addition, neighbors are just unlabeled anonymous individuals who cannot be recognized from one round to the next. In this respect, it is worth noting that most of the reputation subjects assign to partners arises from their average cooperativeness without reference to the chronological set of actions<sup>21</sup>. On the other hand, this is also the case in many e-commerce platforms (e.g., Amazon) where only an average success rate of interactions with external sellers is provided. In this sense, our setup reproduces a real-world situation in which a subject interacts with a partner for the first time, i.e., when first-hand information about the partner is not available.

In the RR treatment each round consisted of the following four stages:

1. Action choice
2. Neighborhood modification
3. Link acceptance decision
4. Feedback on payoffs

In the first stage, players receive information on the cooperation index of their current neighbors and have to select one of two actions,  $A$  or  $B$ , where  $A$  implied "cooperation" and  $B$  implied "defection", the action being the same with all neighbors, as said above. We chose to label actions in a neutral fashion to prevent framing effects<sup>32,33</sup>. In the second stage, participants may decide to unilaterally suppress a link with a neighbor and they are also given the option to offer a link to a new, randomly chosen partner; in both cases, they only know the  $\alpha$  value of the corresponding subject. In the following stage, participants see all link proposals from other players (and their  $\alpha$ ), which they can either accept or reject. After these decision stages a new network is formed, and subjects accumulate their payoff by playing the PD game in pairs with their current neighbors. They are neither informed about their neighbors' payoffs nor about their neighbors' individual current actions. Participants never know the full network topology.

The FR treatment is identical to the RR treatment with the following fundamental difference: Participants never know whether the observed cooperation index  $\alpha$  of their partners is the real one or has been modified. Consequently, in this setup there is an additional stage between stages 1 and 2 of the RR treatment during which



**Figure 1.** True and observable cooperation index  $\alpha$  in the whole population aggregating all treatments in the baseline case (RR, black dots) and the case with fake reputation (FR, blue squares). The observable  $\alpha$  in the FR treatment is represented by the red triangles. Error bars represent standard errors of the mean. The difference between final mean values of true  $\alpha$  is not statistically significant [first repetition,  $P = 0.416$ ; both repetitions,  $P = 0.336$ ]. The difference between final mean values of RR true  $\alpha$  and FR observable  $\alpha$  is statistically significant considering both repetitions [first repetition,  $P = 0.138$ ; both repetitions,  $^*P = 0.019$ ].

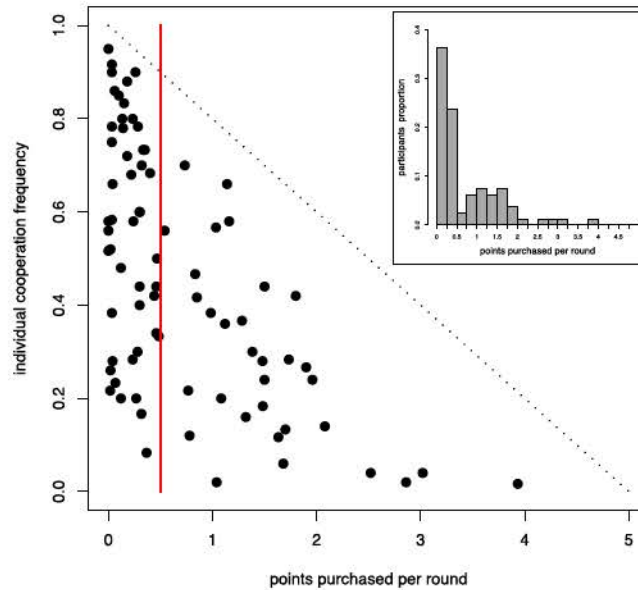
participants may choose to pay a cost in order to modify their  $\alpha$  value. The chosen cost was 4 points for each unity of reputation modified, per round. For example, if a player has currently an  $\alpha$  value of 2 based on her actual last five actions, she can decide to pay 8 points to show an *observable cooperation index* of 4 to the partners. This modification only lasts for the current round. If a player wants to change her observable  $\alpha$  again for the following round she has to pay the cost anew. Apart from that, as in the RR treatment, there is no cost implied if one just wants to show her true cooperation index. Before choosing the above value of four for the cost we performed a preliminary laboratory session in which the cost was set to nine points instead. In that case, we observed that very few players chose to pay that cost to modify their observable  $\alpha$ . Conversely, if the cost is too small then the players would cheat too frequently which would make the cooperation index signal almost useless.

We performed the RR treatment six times where three groups of 20 participants performed the same experiment two times each. The FR treatment was run eight times by four groups each playing two times. Before each new session, we re-initialized the regular lattice by reshuffling the participants who played the same experiment in the same treatment condition for other 30 rounds.

## Results

We now turn to the discussion of our experimental results. First, we look at the behavior of the average cooperation index  $\alpha$  for the baseline case (RR) and for the fake reputation case (FR), see Fig. 1. The time evolution of cooperation in the population, which is noisier, parallels that of  $\alpha$  and it can be found in the SI, Fig. S1.

We now compare the aggregate cooperation frequency results with those obtained in similar recent experimental studies<sup>17–21,23</sup>. However, one must bear in mind that, although the settings are similar in the sense that participants can cut or form links at different rates, the details differ either in link updating frequency, partner accepting rules, information available to the players and, most importantly, the PD payoff matrix values used. In Rand *et al.*<sup>17</sup>, the “fluid dynamic network” treatment is similar to ours, although links to cut and to create are randomly chosen and presented to the players. The information set is also different: the focal player knows the last action of the player at the other end of a random link. In these conditions, Rand *et al.* find that the cooperation frequency stays around 0.6 during all rounds. In Wang *et al.*<sup>18</sup> players update their links at various rates. Information consists in the knowledge of the last five moves of all players. Cooperation stays high at the beginning (more than 0.8) for almost all update frequencies and tends to decay in the final rounds. This behavior is rather expected since this is the only study among those mentioned in which the participants know the exact number of rounds and they are thus eager to defect in the last ones. In Antonioni *et al.*<sup>20</sup> information on the last action of a potential neighbor is costly to participants and it strongly influences the outcome of the experiment. In fact, final cooperation frequencies oscillate between 0.4 and 0.6 for the two values of the cost. On the other hand, when this information is costless cooperation frequency can reach 0.8–0.9. In Cuesta *et al.*<sup>21</sup> the authors investigate how the amount of reputation available influences cooperation in a dynamical environment in which unwanted links can be cut and new ones formed in a manner qualitatively similar to all previously described settings. Reputation is given by the sequence of the last  $m$  actions of any given player where  $m$  can be varied between 0 and 5. The authors find that there is a clear positive correlation between  $m$  and the cooperation level. For  $m = 0$  cooperation quickly decays from an initial 0.5 to 0.2 at the end of the runs. On the other hand when



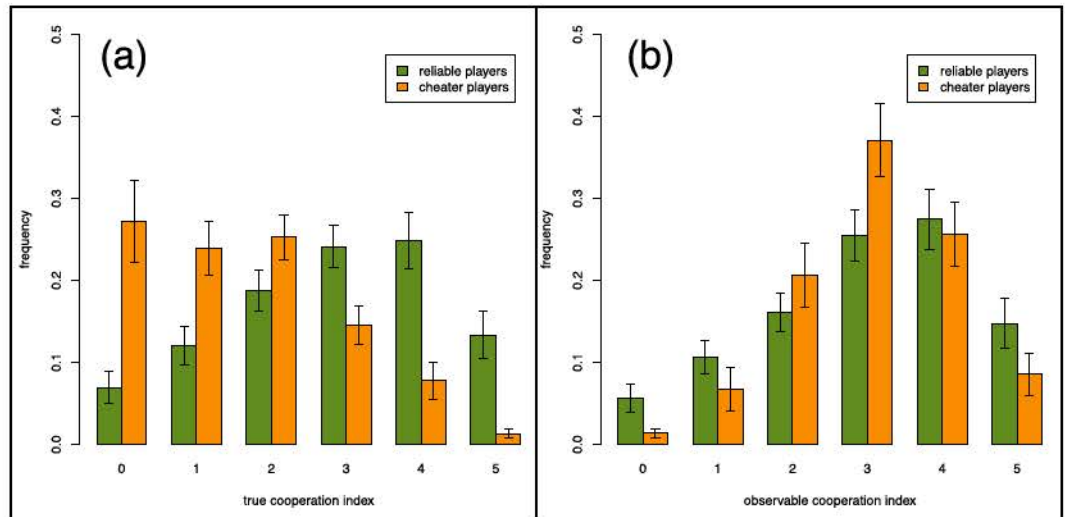
**Figure 2.** Scatterplot of the participants main behavioral features in the FR treatment. The  $x$ -axis value is the average number of points that a given player has paid per round while the  $y$ -axis represents her frequency of cooperation. The red line separates the area containing participants we have called *reliable* (left side) players from the so-called *cheaters* (right side), while the dotted diagonal limits the feasible space a player can be in. Inset: histogram of the proportion of participants who buy a certain amount of points per round on average.

$m > 0$  cooperation is sustained with  $m = 3$  and  $m = 5$  giving statistically indistinguishable results with a roughly constant cooperation level between 0.5 and 0.6. In Gallo and Yan<sup>23</sup> there are four treatments which differ in the amount of information participants have about their partners and about the whole network. In the baseline treatment subjects only know the previous five actions of their direct neighbors, while in the most information-rich environment they know the previous five actions of all players, as well as the topological structure of the current network. The remaining two settings are in between the previous ones. Concerning the level of cooperation, they found that global reputational knowledge is the main determinant for the sustenance of cooperation, which stays at about 0.5–0.6 over the whole period. Knowledge of the structure of the whole network does not help. By contrast, in the setting in which reputational knowledge is only local cooperation stays at about 0.3. Finally, in Fehl *et al.*<sup>19</sup> cooperation reaches high levels around 0.7 but their setting cannot be compared with ours, nor with the above ones because agents there can choose a different action with different neighbors.

With respect to the above-mentioned studies where cooperation is high and remains stable in dynamical networks when information about the partners' strategy is complete, in our case cooperation is maintained but at a lower level (see also Fig. S1 in the SI). We believe that the reason for this difference is to be found in the higher level of uncertainty. Even when  $\alpha$  cannot be faked (RR treatment), the single index that people see being an average and not the true temporal sequence of actions, does not allow cooperative acts to be identified with certainty and participants are left guessing to some extent. In fact, all sessions started with a fraction of cooperators of about 0.6 and this fraction was about 0.5 at the end (see also Fig. S1 in the SI). On the other hand, as shown in<sup>21</sup>, knowledge of the last action of possible partners plays an important part in reputation assignment, going from almost 30% when information comprises the last 3 actions to more than 16% with 5 actions. This missing piece of information may lead subjects to estimate their counterparts' reputation to be lower than what they would do with more information, and therefore to decrease their cooperativeness. Whatever the case, it is important to notice that cooperation based on this kind of easily manipulable reputation system still seems to be fairly high, although our results are not conclusive about the possibility that it will eventually decay. Hence, at least as far as first interactions are concerned, we did not observe a serious hampering of the willingness to cooperate. Other explanations on the observed cooperative behavior are also possible e.g., the influence of the payoff matrix values<sup>34,35</sup> and group sizes<sup>36,37</sup>. Unfortunately, we were not able to run another setting because of time and financial constraints.

Let us now move into between-subject differences in behavior. To that end, in Fig. S3 (see SI) we analyze the average participants' frequency of cooperation in deciles for the RR treatment (black bars) and for the FR treatment (blue bars). Interestingly, it can be seen that in the RR treatment about one third of the participants cooperate between 50% and 60% of the times. Such a peak of cooperation is not observed in the FR treatment where the frequencies tend to be more uniform. In fact, in the FR treatment some participants decided to maintain a lower cooperation frequency and to increase their observable cooperation index paying the cost.

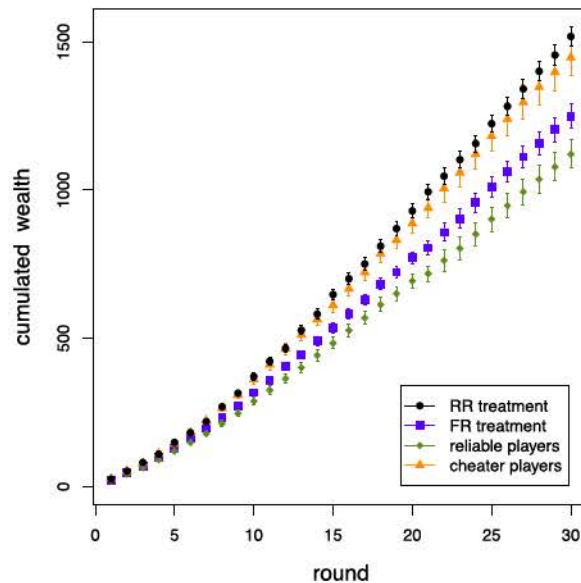
Figure 2 sheds more light on this issue by representing the position of each player in a space where the  $x$ -coordinate is the player's average number of points paid per round and the  $y$ -coordinate is her cooperation frequency for all sessions of the FR treatment. It can be clearly seen that most players cheat only rarely, buying less than half a point per round. Thus we have, somewhat arbitrarily but sensibly, traced the dividing line at this point.



**Figure 3.** Frequency of experimental cooperation indices in the FR treatment separately for cheater and reliable players and for all treatments and rounds. (a) The panel depicts the frequency of the true cooperation index; (b) the panel shows the observable cooperation index. Note that, while reliables behave coherently and have similar  $\alpha$  profiles, cheaters cooperate much less but tend to show an observable cooperation index comparable to that of reliables. The difference in distribution between true cooperation indices is always statistically significant when observed at the beginning and the end of the treatment [first repetition at first round,  $**P = 0.003$ ; both repetitions,  $***P < 0.001$ ; first repetition at last round,  $***P < 0.001$ ; both repetitions,  $***P < 0.001$ ]. The difference in distribution between observable cooperation indices is never statistically significant [1st at first round,  $P = 0.509$ ; both,  $P = 0.640$ ; 1st at last round,  $P = 0.985$ ; both,  $P = 0.388$ ]. The difference in distribution for reliable players is never statistically significant [ $P > 0.9$ ] while for cheater players is always statistically significant [ $***P < 0.001$ ] at the beginning and the end of the treatment. Error bars represent standard errors of the mean.

As will be shown below, this criterion does reflect very well the two main types of behavior in the population. We have dubbed the players that appear in the area to the left of this red line “reliables”. By contrast, the rest of the players, those who buy more than half a point per round in the average, will be called “cheaters”. Among reliables, we observe a heterogeneous behavior: Some are essentially cooperators (top left corner of the scatterplot), some are mostly defectors (bottom left corner), and the rest have a mixed behavior. Cheaters, on the other hand, cooperate less on average. The inset summarizes this information in an aggregate manner by showing the proportion of players that buy a certain amount of points per round. Reliable players are seen to be around 60% of the total. Most cheaters buy between 0.5 and 2.0 points per round, while very few increase their  $\alpha$  by more than two points per round. This suggests that most cheaters tend to stick to an observable cooperation index of about three in order not to trigger link cutting from neighbors. Also, we plot the normalized number of players purchasing points per round and per participant type in Fig. S8 of the SI.

In order to confirm the above hypothesis we plot in Fig. 3 the histograms of the cooperation index in the population for the two categories of subjects. We expect that, if our definition makes sense, reliable players should have very similar true and observable  $\alpha$ , and this is indeed the case. This does not hold for cheaters, who tend to increase their observable  $\alpha$  when the true one is 0 or 1 comparing the histograms in Fig. 3a,b. This is quite understandable given the setting of the experiment and supports our interpretation, namely that the participants’ apparent goal seems to be to show a cooperation index of about three, which guarantees a “fair” behavior on the part of the neighbors who will not be tempted to cut their link to them. In this respect, it is important to note that the general appearance of the histogram of observable  $\alpha$  for cheaters is very similar to that of reliables (Fig. 3b), indicating that cheaters grasp what the acceptable behavior should be. If any, the main difference is that the histogram for cheaters is more peaked around three, i.e., there are fewer cheaters showing a very high  $\alpha$ , in agreement with our intuition that they do not need to look very cooperative. In our experiment we have noticed that participants sever a link about 20% of the time both in the RR and FR treatments. In this respect, Fig. S6 shows that subjects are quite heterogeneous in their link cutting frequency, with a majority of them cutting links less than 20% of the time while others sever their connections much more often (even up to 80% in some cases). From Fig. 3b we can also infer that cheaters use their observable  $\alpha$  to avoid having their links cut off as they appear to be “reliable”. Their observable  $\alpha$  also helps them to be accepted by other players in the link proposal phase. This can also be inferred from Fig. S7, that shows that the higher the  $\alpha$  of an individual, the more likely her acceptance as a new neighbor. Interestingly, in the FR treatment new links have a smaller acceptance rate than in the RR treatment, which is most probably due to the uncertainty about the observable  $\alpha$ . On the other hand, the differences between reliables and cheaters give rise also to noticeable traces on the aggregate behavior: The cooperation level of reliables is considerably larger than that of cheaters (cf. Fig. S1 in the SI), while their combination mimics the true cooperation index of the RR treatment (cf. Fig. S2a in the SI). Again, the observable cooperation index turns out to be almost the same for both kind of players (cf. Fig. S2b in the SI). It is also remarkable that reliables exhibit



**Figure 4.** Cumulated participants' wealth averaged over all sessions for the RR (black dots) and FR (blue squares) treatments. In the latter case we also plot the wealth for reliables (green squares), and cheaters (orange triangles) separately. The cost for reputation modification is taken into account. Error bars represent standard errors of the mean. The difference between final mean values of cumulated wealth for RR and FR treatment is statistically significant considering both repetitions [first repetition,  $P = 0.189$ ; both repetitions,  $*P = 0.023$ ].

a somewhat larger cooperation, as if their honest behavior in terms of reputation would be associated to higher cooperativeness (in a manner not unrelated to the “phenotypes” reported by Peysakhovich *et al.*<sup>38</sup>). We also mention the work by Biziou-van-pol *et al.*<sup>39</sup> who studied the relation between cooperation, altruism, and aversion to telling white lies. White lies are those that increase the benefit of the liar and/or somebody else. Specifically, the authors find that there is a negative correlation between telling white lies that benefit both the other person and the liar and cooperative behavior. This is in line with our findings but one should bear in mind that our subjects tell so called “black lies”, which are those that increase a person's benefit at the expense of another.

So far, we have reported that while the level of cooperation observed in the two treatments is basically the same, the population in the FR treatment shows a clear splitting in two subpopulations, reliables and cheaters. Does this segregation lead to noticeable consequences at the population level? To answer this question, Fig. 4 shows the experimental average cumulated payoff, or social wealth, by treatment and type of player as a function of the round number. First, we note that participants have the best payoff in the RR treatment (black dots). We interpret this result as being a consequence of three factors: a slightly higher cooperation level in the RR treatment (see Fig. 1), the absence of a cost to increase one's reputation, and a slightly higher average degree of the players network in the RR case (cf. SI, Fig. S4). Likewise, the blue squares in Fig. 4 reports the cumulated wealth in the FR treatment whereas the green and orange symbols show, respectively, the cumulated gain for reliables and cheaters taking the cost into account. An interesting result is that cheaters gain more than reliables in spite of paying the cost of cheating. Having a higher reputation allows a cheater to maintain and create more connections to neighbors (cf. SI, Fig. S4) which tends to increase her payoff. Furthermore, a cheater tends to defect more often and thus to earn the maximum payoff  $T$  in many encounters. Thus, although when cheating is possible the total gain is less than in the RR treatment, we see now that it is more profitable to be a cheater in the FR setting. As a result, the inequality in our “society” increases: The fact that cheaters earn a higher payoff leads to a Gini coefficient of 0.370 in the FR treatment, to be compared to a value of 0.271 in the RR treatment, both indices being calculated on the cumulated wealth of participants at the end of the experiment.

For completeness, we also show the time evolution of the participants' payoff in the SI, Fig. S5. Now, comparing the cumulated gain of cheaters in the FR treatment with the cumulated gain in the RR treatment shows that there is little difference but the payoff is slightly larger in the latter. Thus, although it pays to be a cheater when faking one's reputation is allowed, if nobody is allowed to do it the social wealth is higher, at least as far as this experiment is concerned.

## Discussion

In summary, we have designed and carried out an experiment in order to test the effects of uncertainty about the reputation of possible partners in the frame of the Prisoner's Dilemma game. Our experiment provides us with enough evidence to support several important conclusions. To begin with, the aggregate cooperation level of the population does not change when reputations can be faked. Interestingly, this is in agreement with the only game-theoretical work we know of in this context: Röhl *et al.*<sup>40</sup> showed, in an evolutionary public goods game, that fake reputation does not harm cooperation under some conditions in well-mixed populations. We

must stress that their approach is quite different from ours, in particular because every individual interacts with every other one and cannot modify this interaction neighborhood. On the other hand, Röhl *et al.* introduce a probability to be discovered and punished which is reminiscent of the link-cutting stage of our experiment: Note that our subjects cut links without knowing for sure that the corresponding individual is a cheater, which is not unrelated to probabilistic discovery. Therefore, in spite of the differences, the fact that our observation aligns with the predictions in<sup>40</sup> is certainly suggestive. Another interesting line of research related to the findings we are reporting here is that of image scoring in evolutionary games<sup>24,41</sup>, based on the idea that helping someone increases one's image score, whereas refusing to help reduces it. This is clearly similar to the notion of reputation, except that when it is kept by each individual separately it becomes private instead of public as we consider here. In this regard, it has been recently shown<sup>27</sup> that when there is information on group scoring only a few images are needed to sustain reputation. This points to the possibility of preventing the problems of faking reputation by externally providing some manner of (truthful) group information.

Nevertheless, our results go further than this as they point to a splitting of the experimental population in two different types of individuals: reliables, who cheat very little if at all, and cheaters, who are willing to fake their reputation almost at every interaction. This cannot be noticed by looking at the cooperation level and only an analysis of the within-subject variability allows to uncover this effect. In spite of the fact that cheaters have to pay some cost to modify their cooperation index, they still end up making more profit than reliables, as they manage to exhibit an intermediate reputation that makes them less likely targets for link cutting and more likely to be accepted as new partners. The similarity between the histogram of the cooperation index in both treatments is striking and proves that cheaters have a correct intuition about what is the optimum level of reputation modification to do well. In addition, it turns out that the average earnings of cheaters are very similar to the case in which reputation is truthful, which implies that reliables are in practice paying the cost of the cheaters' efforts to disguise their bad behavior. On the other hand, reliables are honest not only with respect to their reputation, but also about being even more cooperative than the general population in the baseline treatment. As a result of the combination of reliable and cheating behaviors, inequality increases: the Gini coefficient of the RR treatment increases by more than a 30% in the FR treatment; for a comparison which only has illustrative value, these would be values similar, respectively, to Finland in 2008 and Tanzania in 2007<sup>42</sup>. We therefore conclude that, in our experiment, even if the level of cooperation in interactions is basically the same as when reputation is truthful, the features that emerge from the possibility of faking reputations are the splitting into exploiters and exploited, and a larger degree of inequality, both highly undesirable. While this is a first step in the experimental analysis of the problem that cannot be easily generalized to more complex social environments, it is clear that the conclusions are still very relevant and would justify further, intensive research along these lines, in order to inform policy makers' efforts to ensure fair and transparent trade.

## Methods

The use of human subjects in this experiment has been approved by the Ethics Committee of the University of Lausanne and our methods were carried out in accordance with the approved guidelines. Participants signed an informed consent describing the nature of the experiment before they entered into the laboratory. We conducted a total of seven experimental sessions in November 2014. Participants were recruited from the pool of undergraduate students from all disciplines of the University of Lausanne and the Ecole Polytechnique Fédérale de Lausanne using ORSEE<sup>43</sup>. Subject-subject anonymity was granted at all stages, and the experiment was computerized using the z-Tree environment<sup>44</sup>. Before making decisions, participants read detailed instructions and responded to a set of control questions in order to insure common understanding of the game and the computation of payoffs. A translation of these instructions from the original French is provided in Section S1 of the SI. Each session lasted one and a half hours and included 20 participants, where a total of 140 subjects, 48 women and 92 men, took part in the experiment. Participants were randomly assigned to the RR treatment (60 subjects) and to the FR treatment (80 subjects). Subjects observable demographic variables did not qualitatively differ across treatments. Participants received a show-up fee of 10 CHF (about 10.50\$), and their final score in points was converted at an exchange rate of 1 CHF = 120 points. The average payoff per subject was 27.35 CHF (about 29\$). All statistical difference significances of mean values have been obtained performing unequal variances *t*-test analysis. All statistical difference in distribution have been obtained performing two-sample Kolmogorov-Smirnov test analysis. We considered the fifth one as the first observable and independent round to compare cooperation index distributions. Since each group performed two repetitions of the assigned treatment we considered two different statistical approaches. The first one takes into account independent observations, thus considering only the first repetition of the treatment, while the second analysis assumes both repetitions of the treatment as two independent observations.

## References

1. Rainie, L. & Wellman, B. *Networked. The New Social Operating System* (MIT Press, Cambridge, MA, 2012).
2. Kendall, L. "Community and the internet". In *The Handbook of Internet Studies* (eds Consalvo M. *et al.*) 309–325 (Wiley-Blackwell, 2011).
3. van Dijck, J. *The Culture of Connectivity: A Critical History of Social Media* (Oxford University Press, Oxford, UK, 2013).
4. Stephany, A. *The Business of Sharing: Making it in the New Sharing Economy* (Palgrave Macmillan, London, UK, 2015).
5. Streitfeld, D. Give yourself 5 stars? Online, it might cost you. *The New York Times* (2013) Available at: <http://goo.gl/0Tuz7U> (Accessed: 5th May 2016).
6. Conner, C. The dark side of reputation management: How it affects your business. *Forbes* (2013) Available at: <http://goo.gl/NZ88Yz> (Accessed: 5th May 2016).
7. Merrit, J. *Fighting fake reviews: Removal, response and your reputation*. (2014) Available at: <https://goo.gl/iFqCcc> (Accessed: 5th May 2016).
8. Rapoport, A. & Chammah, A. M. *Prisoner's Dilemma* (University of Michigan Press, Ann Arbor, 1965).



9. Axelrod, R. *The Evolution of Cooperation* (Basic Books, Inc., New York, 1984).
10. Nowak, M. A. Five rules for the evolution of cooperation. *Science* **314**, 1560–1563 (2006).
11. Fletcher, J. A. & Doebeli, M. A simple and general explanation for the evolution of altruism. *Proc. R. Soc. B* **276**, 13–19 (2009). dilemma experiments: Conditional cooperation and payoff irrelevance. *Sci. Rep.* **4**, 4615 (2014).
12. Skyrms, B. & Pemantle, R. A dynamic model for social network formation. *Proc. Natl. Acad. Sci. USA* **97**, 9340–9346 (2000).
13. Eguiluz, V. M., Zimmermann, M. G., Cela-Conde, C. J. & San Miguel, M. Cooperation and the emergence of role differentiation in the dynamics of social networks. *Am. J. Sociol.* **110**, 977–1008 (2005).
14. Santos, F. C., Pacheco J. M. & Lenaerts, T. Cooperation prevails when individuals adjust their social ties. *PLoS Comput. Biol.* **2**, 1284–1291 (2006).
15. Perc, M. & Szolnoki, A. Coevolutionary games - A mini review. *Biosystems* **99**, 109–125 (2010).
16. Perc, M., Gómez-Gardeñes, J., Szolnoki, A., Floría, L. M. & Moreno, Y. Evolutionary dynamics of group interactions on structured populations: A review. *J. R. Soc. Interface* **10**, 20120997 (2013).
17. Rand, D. G., Arbesman, S. & Christakis, N. A. Dynamic social networks promote cooperation in experiments with humans. *Proc. Natl. Acad. Sci. USA* **108**, 19193–19198 (2011).
18. Wang, J., Suri, S. & Watts, D. J. Cooperation and assortativity with dynamic partner updating. *Proc. Natl. Acad. Sci. USA* **109**, 14363–14368 (2012).
19. Fehel, K., van der Post, D. J. & Semmann, D. J. Co-evolution of behavior and social network structure promotes human cooperation. *Ecol. Lett.* **14**, 546–551 (2011).
20. Antonioni, A., Cacaull, M. P., Lalive, R. & Tomassini, M. Know thy neighbor: Costly information can hurt cooperation in dynamic networks. *PLOS ONE* **9**, e110788 (2014).
21. Cuesta, J. A., Gracia-Lázaro, C., Ferrer, A., Moreno, Y. & Sánchez, A. Reputation drives cooperative behaviour and network formation in human groups. *Sci. Rep.* **5**, 7843 (2014).
22. Gracia-Lázaro, C. *et al.* Heterogeneous networks do not promote cooperation when humans play a Prisoner's Dilemma. *Proc. Natl. Acad. Sci. USA* **109**, 12922–12926 (2012).
23. Gallo, E. & Yan, C. The effects of reputational and social knowledge on cooperation. *Proc. Natl. Acad. Sci. USA* **112**, 3647–3652 (2015).
24. Wedekind, C. & Milinski, M. Cooperation through image scoring in humans. *Science* **288**, 850–852 (2000).
25. Milinski, M., Semmann, D. & Krambeck, H. J. Reputation helps solve 'the tragedy of the commons'. *Nature* **415**, 424–426 (2002).
26. Sommerfeld, R. D., Krambeck, H. J. & Milinski, M. Multiple gossip statements and their effect on reputation and trustworthiness. *Proc. Roy. Soc. B* **275**, 2529–2536 (2008).
27. Nax, H. H., Perc, M., Szolnoki, A. & Helbing, D. Stability of cooperation under image scoring in group interactions. *Scientific Reports* **5**, 12145 (2015).
28. Nowak, M. A. & Sigmund, K. Evolution of indirect reciprocity by image scoring. *Nature* **393**, 573–577 (1998).
29. Nowak, M. A. & Sigmund, K. Evolution of indirect reciprocity. *Nature* **437**, 1291–1298 (2005).
30. Camerer, C. F. *Behavioral Game Theory: Experiments in Strategic Interaction* (Princeton University Press, Princeton, 2003).
31. Montague, D. A. *Essentials of Online Payment Security and Fraud Prevention* (John Wiley & Sons, Hoboken, NJ, 2010).
32. Liberman, V., Samuels, S. M. & Ross, L. The name of the game: predictive power of reputations versus situational labels in determining Prisoner's Dilemma game moves. *Pers. Soc. Psychol. B.* **30**, 1175–1185 (2004).
33. Ellingsen, T., Johannesson, M., Mollerstrom, J. & Munkhammar, S. Social framing effects: preferences or beliefs? *Games Econ. Behav.* **76**, 117–130 (2012).
34. Gunthorsdottir, A., Houser, D. & McCabe, K. Dispositions, history and contributions in public goods experiments. *J. Econ. Behav. Organ.* **62**, 304–315 (2007).
35. Capraro, V., Jordan, J. J. & Rand, D. G. Heuristics guide the implementation of social preferences in one-shot Prisoner's Dilemma games. *Sci. Rep.* **4**, 6790 (2014).
36. Szolnoki, A. & Perc, M. Impact of critical mass on the evolution of cooperation in spatial public goods games. *Phys. Rev. E* **81**, 057101 (2010).
37. Barcelo, H. & Capraro, V. Group size effect on cooperation in one-shot social dilemmas. *Sci. Rep.* **5**, 7937 (2015).
38. Peysakhovich, A., Nowak, M. A. & Rand, D. G. Humans display a cooperative phenotype that is domain general and temporally stable. *Nat. Commun.* **5**, 4939 (2014).
39. Biziou-van-Pol, L., Haenen, J., Novaro, A., Occhipinti-Liberman, A. & Capraro, V. Does telling white lies signal pro-social preferences? *Judgm. Decis. Mak.* **10**, 538–548 (2015).
40. Röhl, T., Röhl, C., Schuster, H. G. & Traulsen, A. Impact of fraud on the mean-field dynamics of cooperative social systems. *Phys. Rev. E* **76**, 026114 (2007).
41. Milinski, M., Semmann, D., Bakker, T. C. M. & Krambeck, H.-J. Cooperation through indirect reciprocity: image scoring or standing strategy? *Proc. Roy. Soc. B* **268**, 2495 (2001).
42. CIA. *The World Factbook 2013-14* (Central Intelligence Agency, Washington, DC, 2013).
43. Greiner, B. Subject pool recruitment procedures: Organizing experiments with ORSEE. *J. Econ. Sci. Ass.* **1**, 114–125 (2015).
44. Fischbacher, U. z-Tree: Zürich toolbox for ready-made economic experiments. *Exp. Econ.* **10**, 171–178 (2007).

## Acknowledgements

A. A. gratefully acknowledges financial support by the Swiss National Science Foundation (under grants no. 200020-143224, CR13I1-138032 and P2LAP1-161864) and by the Rectors' Conference of the Swiss Universities (under grant no. 26058983). All authors acknowledge financial support to carry out the experiments by the Faculty of Business and Economics of the University of Lausanne and the fundamental support by Prof. Rafael Lalive. This work has been supported in part by the European Commission through FET Open RIA 662725 (IBSEN) and by the Ministerio de Economía y Competitividad (Spain) under grant FIS2015-64349-P (VARIANCE).

## Author Contributions

All authors conceived the experimental setting. A.A. ran laboratory experiments and performed the data analysis. All authors discussed the results, drew conclusions and wrote the manuscript.

## Additional Information

Supplementary information accompanies this paper at <http://www.nature.com/srep>

Competing financial interests: The authors declare no competing financial interests.

**How to cite this article:** Antonioni, A. *et al.* Cooperation Survives and Cheating Pays in a Dynamic Network Structure with Unreliable Reputation. *Sci. Rep.* **6**, 27160; doi: 10.1038/srep27160 (2016).



This work is licensed under a Creative Commons Attribution 4.0 International License. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in the credit line; if the material is not included under the Creative Commons license, users will need to obtain permission from the license holder to reproduce the material. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>

# Cooperation Survives and Cheating Pays in a Dynamic Network Structure with Unreliable Reputation

## Supplementary Information

A. Antonioni, A. Sánchez, and M. Tomassini

The Supplementary Information (SI) provides translation from French of the exact instructions form that participants received for the FR treatment and corresponding trial questions, see Section 1. Instructions for the baseline RR treatment are identical without the *apparent profile* stage.

Additional results of experimental data that are not shown in the main text are presented in Section 2.

# 1 Instructions (translated from French)

Welcome to this experiment !

You are going to take decisions that will affect your payoffs as well as the payoffs of the other participants.

Although all payoffs are expressed in number of points, these points will be transformed into money at the end of the experiment according to the following exchange rate:

$$120 \text{ pts.} = 1.- \text{ CHF}$$

During the experiment it is **strictly forbidden to talk to other participants**. If you have a question, please ask the assistants. If you don't comply with these rules, we will regrettably be obliged to exclude you from the experiment.

In this experiment every participant will interact with his/her "neighbors" in a network constituted by all the participants in the room. At the beginning of the experiment everybody will have four neighbors but this number may change during the experiment as explained below.

During the experiment you'll only see your immediate neighbors, i.e. only the participants with whom you are directly linked; you will not see what happens in the rest of the network (for instance, you won't be able to see what the neighbors of your direct neighbors are doing).

## What is it all about?

There will be a number of rounds that is comprised between 20 and 50 but the exact number will not be made explicit. Each round consists of four steps:

1. Decide your action
2. Decide your apparent profile
3. Modify your neighborhood
4. Accept or refuse new links

## 1. Decide your action

In this first stage you will have to choose an “action” among the two following options:

**A or B**

Like you, your neighbors will also have to take the same decision about the action they will choose. The chosen action is unique: this means that the same action will be used when interacting with all of your neighbors (you can't use different actions with different neighbors). Your payoff for the current round is computed as a function of your current action and the current actions of your neighbors.

Now we explain the payoffs for each possible combination of your action and one's of your neighbors action:

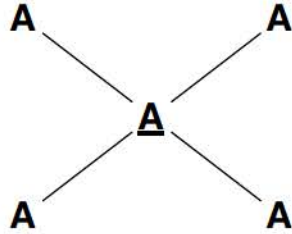
- You choose **A** and your neighbor chooses **A**:  
you gain **7 points** and your neighbor also gets **7 points**.
- You choose **A** and your neighbor chooses **B**:  
you get **0 points** and your neighbor gets **10 points**.
- You choose **B** and your neighbor chooses **A**:  
you gain **10 points** and your neighbor gains **0 points**.
- You choose **B** and your neighbor chooses **B**:  
you get **0 points** and your neighbor gets **0 points**.

Your final accumulated payoff in each round is computed as the sum of the points gained in each interaction with each of your current neighbors.

However, the relevant neighbors for the payoff computation are those to whom you are directly linked **at the end of each round**. As you shall see below, you will be allowed to modify your neighborhood in the following stages before the round ends.

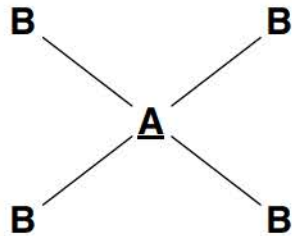
The following examples illustrate the computation of your payoff at the end of a round. Note that in the examples you are the central player, with your action underlined, and you have **4 neighbors** but during the experiment you could have a different number of neighbors.

**Example 1** : Your action is **A**, the action of all your neighbors is also **A**.



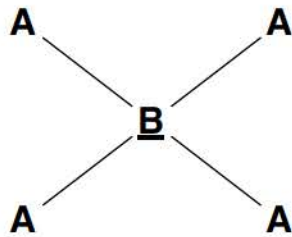
Your payoff is :  $7 + 7 + 7 + 7 = 28$  points.

**Example 2** : Your action is **A**, the action of all your neighbors is **B**.



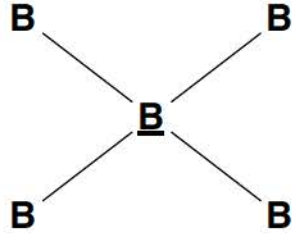
Your payoff is :  $0 + 0 + 0 + 0 = 0$  points.

**Example 3** : Your action is **B**, the action of all your neighbors is **A**.



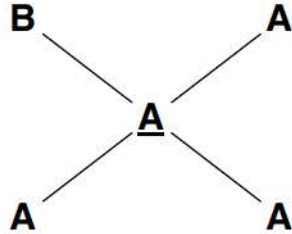
Your payoff is :  $10 + 10 + 10 + 10 = 40$  points.

**Example 4** : Your action is **B**, the action of all your neighbors is also **B**.



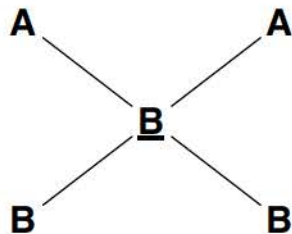
Your payoff is :  $0 + 0 + 0 + 0 = 0$  points.

**Example 5** : Your action is **A**, the action of three of your neighbors is **A** while the fourth one does **B**.



Your payoff is :  $0 + 7 + 7 + 7 = 21$  points.

**Example 6** : Your action is **B**, two of your neighbors do **A** while the other two do **B**.



Your payoff is :  $10 + 10 + 0 + 0 = 20$  points.

## The Profile

The *profile* of a participant is computed from the sequence of her/his last five actions. It is defined as the number of **A**s in the sequence and, by consequence, it belongs to the integer interval  $[0, 5]$ .

For example, if the past actions sequence is:

A	B	B	A	B
---	---	---	---	---

the corresponding profile value is **2**.

If the past actions sequence is:

B	A	A	A	B
---	---	---	---	---

then the corresponding profile value is **3**.

The above definition refers to the profile called *real*.

You have the option of choosing an *apparent* profile different from the real one at each round, according to the procedure explained in the next section of this document. This apparent profile is the one that your neighbors will see. The real profile is private and will not be communicated to your neighbors unless your apparent profile is the real one, i.e. you choose not to alter it in the current round.

On the other hand, note that the current actions of your neighbors will never be known to you, only their apparent profiles.

## IMPORTANT

Note that at the beginning of the experiment all participants will start with the same initial sequence of actions as follows:

A	B	A	B	A
---	---	---	---	---

which amounts to an initial profile value of **3** for all.

During the experiment the profile will be denoted by "P".



In the first stage of each round you must choose the action that you want to use during the current round with the help of the following screenshot:

The screenshot shows a game interface with a timer in the top right corner indicating 20 seconds. At the top, there are five buttons labeled 'A', 'B', 'A', 'B', and 'A'. Below the first and last buttons are labels: 'action la plus ancienne' and 'action la plus récente'. The main area contains two panels. The left panel displays 'Votre profil réel: 3' and 'Votre profil apparent: 3', followed by a list of neighbors with profiles P=0 to P=5 and their counts. The right panel shows 'Votre action pendant ce tour:' with radio buttons for 'A' and 'B', and an 'OK' button.

Temps [sec] 20

A B A B A

action la plus ancienne action la plus récente

Votre profil réel: 3  
Votre profil apparent: 3

Vous avez 4 voisin(s) avec le(s) profile(s) suivant(s):

Profil	Voisins avec ce profil
P = 0:	0
P = 1:	0
P = 2:	0
P = 3:	4
P = 4:	0
P = 5:	0

Votre gain cumulé: 0

Votre action pendant ce tour:

A  
 B

OK

Choose your action for this round and click the "OK" button.

## 2. Decide your apparent profile

In this stage you must decide what to do with your profile for this round.

Important Note: each point added or subtracted to your real profile will cost you 4 points. For example:

- if your profile real value is **2** and you would like to increase it to an apparent profile of **4**, you will have to pay  $2 \times 4 = 8$  points.
- if your real profile is **3** and you want to have an apparent profile of **2** then you must pay  $1 \times 4 = 4$  points.

If you decide to keep your apparent profile as your real profile it will cost you **0** points as this is the default value. If you would like to change your apparent profile but you cannot pay for the corresponding points your apparent profile will be the same as your real profile for the round. Also note that during the first round you cannot change your apparent profile since all participants start with **0 points**.

You will choose your apparent profile for a give round through the following screenshot:

Temps (sec): 15

B A B A A

action la plus ancienne      votre action actuelle

Votre profil réel: 3  
Votre profil apparent: 3

Vous avez 4 voisin(s) avec le(s) profil(s) suivant(s):

Profil	Voisins avec ce profil
P = 0:	0
P = 1:	0
P = 2:	0
P = 3:	4
P = 4:	0
P = 5:	0

Votre gain cumulé: 0

Choisissez votre profil apparent pour ce tour:  
(le choix par défaut à coût zéro est votre profil réel)

0  
 1  
 2  
 3  
 4  
 5

OK

Choose your apparent profile for this round and click the "OK" button. If you want to keep your real profile as apparent profile just click directly the "OK" button.

### 3. Modify your neighborhood

In this stage you must decide whether you want to modify your neighborhood. You'll see in the screen the number of current neighbors you have as well as the values of their apparent profiles. Next, you must decide the following:

- decide whether you want to cut a link to a neighbor with a certain current apparent profile value. An arbitrary link to such a neighbor will be cut automatically and unconditionally.
- decide whether you want to propose a link to a randomly chosen participant who is not yet one of your direct neighbors. The system will made her/his apparent profile known to you. This link will only be created if the partner will accept it in the next stage. In some cases, it will not be possible to propose you creating a new link.

The possible modification of your neighborhood will be made through the following screenshot:

The screenshot shows a game interface with a timer in the top right corner indicating 21 seconds. At the top, there are five buttons labeled 'B', 'A', 'B', 'A', 'A'. Below these, the text 'action la plus ancienne' is on the left and 'votre action actuelle' is on the right. The main area is divided into two panels. The left panel displays the player's current status: 'Votre profil réel: 3', 'Votre profil apparent: 3', and 'Vous avez 4 voisin(s) avec le(s) profil(s) suivant(s):'. Below this is a table of neighbor counts for profiles P=0 to P=5. The right panel contains two decision boxes: 'Voulez-vous couper un lien?' with radio buttons for 'Non' and 'Oui, avec un P=0' through 'P=5', and 'Voulez-vous proposer un lien avec un P=3?' with radio buttons for 'Non' and 'Oui'. An 'OK' button is located at the bottom right of the interface.

Profil	Voisins avec ce profil
P = 0	0
P = 1	0
P = 2	0
P = 3	4
P = 4	0
P = 5	0

Votre gain cumulé: 0

Choose whether you want to cut a link and/or whether you want to propose a new link, then click the "OK" button.

#### 4. Accept or refuse new links

Finally, in this stage you must decide which links, among those proposed by other participants, you are going to accept after knowing the apparent profile of the potential partner. Note that a link will only be created if you accept it at this stage. The proposing participant does already know your apparent profile and decided to propose you creating the link.

You will complete this stage by using the following screenshot:

Temps [sec]: 4

action la plus ancienne: B A B A A

action actuelle: A

Votre profil réel:	3
Votre profil apparent:	3
Vous avez 4 voisin(s) avec le(s) profil(s) suivant(s):	
Profil	Voisins avec ce profil
P = 0:	0
P = 1:	0
P = 2:	0
P = 3:	4
P = 4:	0
P = 5:	0
Votre gain cumulé:	0

Voulez-vous accepter un lien avec un P=3?

Non  
 Oui

Voulez-vous accepter un lien avec un P=37?

Non  
 Oui

OK

For each link proposal, choose whether you want to accept it then click the "OK" button.

## End of round

According to the decisions of all participants, a new network will be created. Your payoff for the round will be computed as a function of your own action and of the actions of all your current neighbors according to the procedure described in the section "Decide your action". From this payoff value it will be subtracted the cost of changing your apparent profile if you decided to do so in this round.

Temps [sec] 25

action la plus ancienne: B A B A A

action actuelle: A

Voire profil réel: 3  
Voire profil apparet: 3

Vous avez 4 voisin(s) avec le(s) profil(s) suivant(s):

Profil	Voisins avec ce profil
P = 0:	0
P = 1:	0
P = 2:	0
P = 3:	4
P = 4:	0
P = 5:	0

Voire gain cumulé: 28

Voire gain dans ce tour: 28  
Voire dépense dans ce tour: 0  
Voire profit dans ce tour: 28

OK

Click the "OK" button to go to the following round.

## Did you understand the explanations ?

Before starting the actual experiment we would like to be sure that you and everybody else has correctly understood the decisions that you'll have to make. To this end, please answer the questions that will appear on your screen. When you are done with a question click the "OK" button at the bottom of the screen.

## Trial questions

Participants answer the following trial questions after reading the instructions. Wrong answers were not accepted by the software in order to continue to the following question. Participants were able to begin the experiment only answering correctly to all questions.

1. How many neighbors do you have at the beginning of the experiment?  
*Correct answer:* 4 neighbors.
2. Suppose that at the end of a round you have 3 neighbors. Your action is **A** and all your neighbors have chosen **A**.  
Which is your gain at the end of this round?  
*Correct answer:* 21 points.
3. Suppose that at the end of a round you have 3 neighbors. Your action is **B** and all your neighbors have chosen **A**.  
Which is your gain at the end of this round?  
*Correct answer:* 30 points.
4. Suppose that at the end of a round you have 5 neighbors. Your action is **A** and three of your neighbors have chosen **B** while the other two of them have chosen **A**.  
Which is your gain at the end of this round?  
*Correct answer:* 14 points.
5. Suppose that at the end of a round you have 5 neighbors. Your action is **B** and three of your neighbors have chosen **B** while the other two of them have chosen **A**.  
Which is your gain at the end of this round?  
*Correct answer:* 20 points.
6. If you decide to **cut off** a link, is the link automatically cut off or is it necessary that also your neighbor accept this decision?  
*Correct answer:* The link is automatically cut off.
7. If you decide to **create** a link, is the link automatically created or is it necessary that also your neighbor accept this decision?  
*Correct answer:* The link is created only if your neighbor accepts it.
8. Suppose that your *real* profile is 3 and you want to have an *apparent* profile of 1. How many points do you need to pay?  
*Correct answer:* 8 points.
9. Suppose that your *real* profile is 0 and you want to have an *apparent* profile of 3. How many points do you need to pay?  
*Correct answer:* 12 points.

## 2 Additional results

The following Fig. S1 shows the average cooperation level for the RR and FR treatments. Moreover, we also show the average cooperative acts of reliable and cheater players, respectively. Figure S2 shows the average cooperation indices (true  $\alpha$ , Fig. S2a; observable  $\alpha$ , Fig. S2b) for the RR and FR treatments, and for reliable and cheater players.

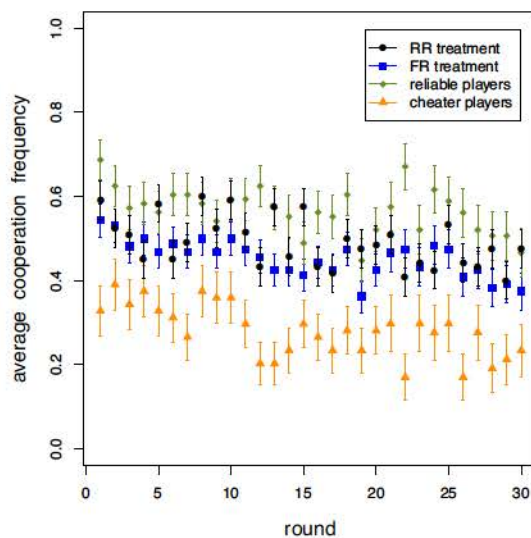


Figure S1: Average cooperation level in the experiment for RR and FR treatments and for reliable and cheater players in the latter case. Error bars represent standard errors of the mean. The difference between final mean values of cooperation level for reliable and cheater players is statistically significant [first repetition,  $P^* = 0.046$ ; both repetitions,  $P^{**} = 0.002$ ]. Other differences are not statistically significant.

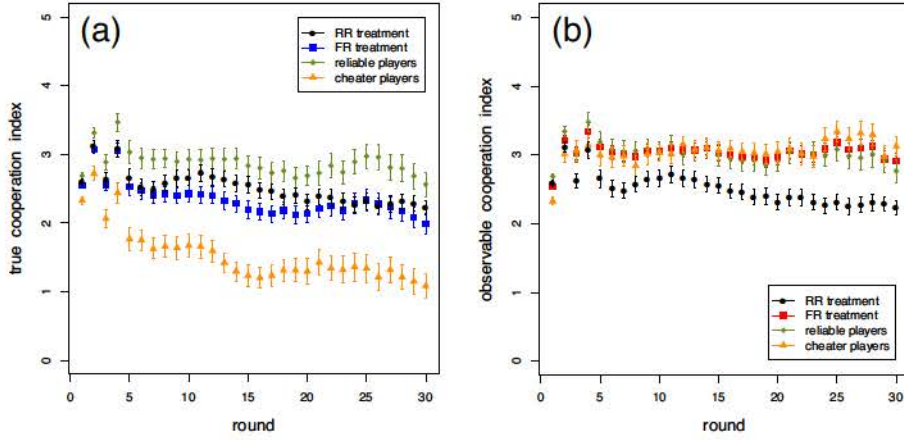


Figure S2: (a): time evolution of the true cooperation index in the whole population after aggregating all the treatments in the baseline case (RR, black dots), the fake reputation treatment (FR, blue squares), and for reliable and cheater players (green squares and orange triangles). (b): time evolution for the observable cooperation index for the same cases. Error bars represent standard errors of the mean. The difference between final mean values of true cooperation index for reliable and cheater players is statistically significant [first repetition,  $P^* = 0.036$ ; both repetitions,  $P^{***} < 0.001$ ]. The difference between final mean values of observable cooperation index for RR and FR treatment is statistically significant considering both repetitions [first repetition,  $P = 0.138$ ; both repetitions,  $P^* = 0.019$ ]. Other differences are not statistically significant.

In Fig. S3 we show the average participants' frequency of cooperation in deciles for the RR treatment (black bars) and for the FR treatment (blue bars). Interestingly, it can be seen that in the RR treatment about one third of the participants cooperate between 50% and 60% of the times. Such a peak of cooperation is not observed in the FR treatment where the frequencies tend to be more uniform. In fact, in the FR treatment some participants decided to maintain a lower cooperation frequency and to increase their observable cooperation index paying the cost.

The following Fig. S4 depicts the evolution of the average degree in the dynamic networks in RR and FR treatments and includes reliable and cheater players. The network mean degree tends to increase but it stabilizes and the graph never becomes a complete one. An interesting effect is that cheaters tend to attract more neighbors in the network since their observable cooperation index is usually higher than the one of reliable players. We note that, for technical reasons, two groups of the FR treatment only performed 20 additional rounds in the second session of the same treatment. All averages have been computed taking this into account.



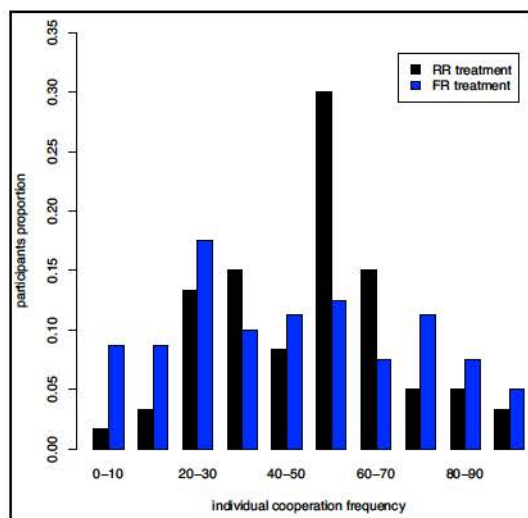


Figure S3: Proportion of participants ( $y$ -axis) who have had a given frequency of cooperative acts during the whole experiment ( $x$ -axis), cumulated over all sessions and grouped in deciles, for the RR and FR treatment.

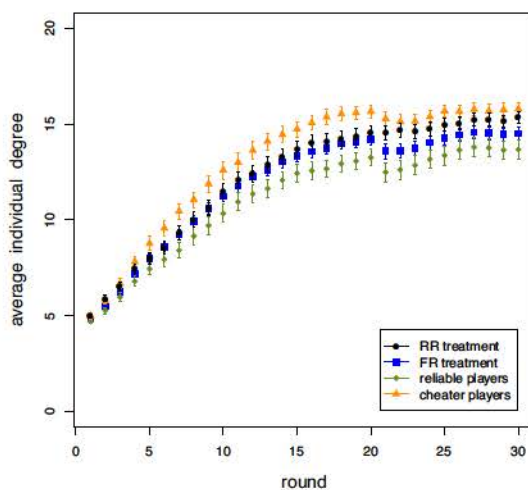


Figure S4: Evolution of the average degree for the RR and FR treatments (all sessions cumulated) and for reliable and cheater players for all sessions in the FR treatment. Error bars represent the standard error of the mean. The difference between final mean values of individual degree for reliable and cheater players is statistically significant [first repetition,  $P^* = 0.046$ ; both repetitions,  $P^{**} = 0.010$ ]. Other differences are not statistically significant.

Figure S5 reports the time evolution of the average payoff and it shows that participants gain more in the RR treatment than in the FR treatment at all times. Conversely, we see that it pays to be a cheater with respect to a reliable player in the FR treatment.

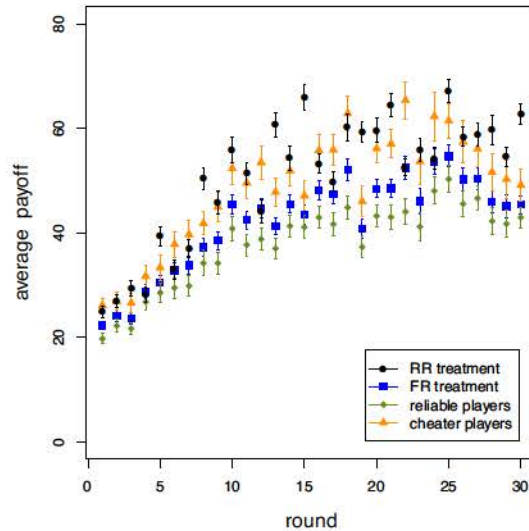


Figure S5: Average payoff per round in the baseline (RR) and in the fake reputation (FR) treatment respectively. We also plot average payoff per round in the FR treatment for reliable and cheater players respectively. Cost is taken into account in these figures. Error bars represent the standard error of the mean. The difference between final mean values of individual payoff for RR and FR treatments is statistically significant considering both approaches [first repetition,  $P^* = 0.072$ ; both repetitions,  $P^* = 0.013$ ]. Other differences are not statistically significant.

Subsequently, Fig. S6 summarizes the main features of individual players by plotting them in a scatterplot given by their link cutting frequency and their cooperativeness in the RR and FR treatment, indicating whether they are reliable or cheaters in the FR treatment. Finally, Fig. S7 depicts the number of accepted proposals in the first stage of the experiment, namely when our experimental setup offers subjects a new link with a randomly chosen player, for the RR and FR treatments.

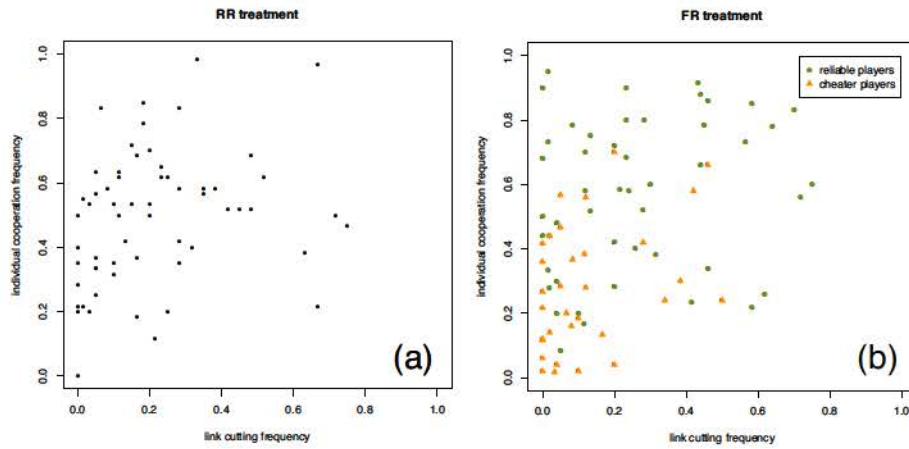


Figure S6: Scatterplots of the players as described by their link cutting frequency ( $x$ -axis) and their cooperation frequency ( $y$ -axis). Each point corresponds to an individual player. (a) RR treatment. (b) FR treatment. Colors in the FR treatment indicate whether a player is a reliable (green dots) or a cheater (orange triangles).

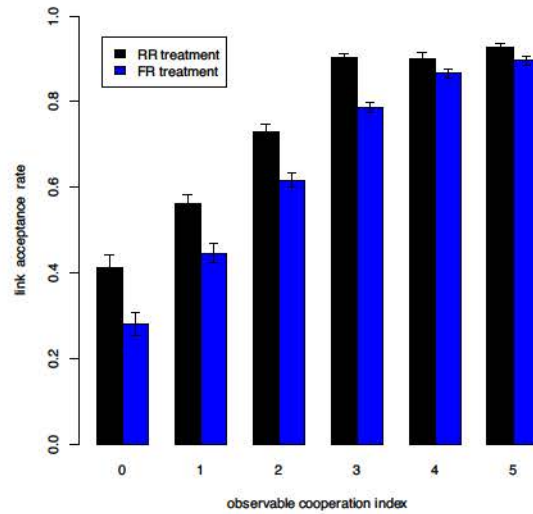


Figure S7: Link acceptance rate at the first stage of the modifications to the network in the experimental setup, i.e., fraction of randomly proposed new neighbors that are accepted as a function of their observable cooperation index. Black bars correspond to the RR treatment, blue bars to the FR treatment. Error bars represent standard error of the mean.

In Fig. S8 we show the normalized number of participants who purchase points (at least one) per round and per participant type in the first repetition of FR treatments. We can see that the vast majority of cheater players fake their cooperation index almost at every interaction, with a particular increase during last rounds. On the other hand, reliable players, by definition, purchased less reputational points and they consequently have a smaller rate.

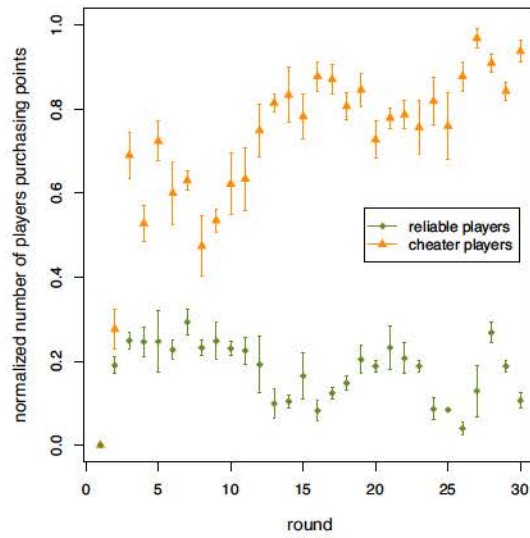


Figure S8: Average normalized number of players who purchase points (at least one) per round and per participant type considering only the first repetition of FR treatments. Error bars represent standard error of the mean.