

Prosocial preferences do not explain human cooperation in public-goods games

Maxwell N. Burton-Chellew^{a,b} and Stuart A. West^{b,1}

^aNuffield College, University of Oxford, Oxford OX1 1NF, United Kingdom; and ^bDepartment of Zoology, University of Oxford, Oxford OX1 3PS, United Kingdom

Edited* by Raghavendra Gadagkar, Indian Institute of Science, Bangalore, India, and approved November 9, 2012 (received for review June 27, 2012)

It has become an accepted paradigm that humans have “prosocial preferences” that lead to higher levels of cooperation than those that would maximize their personal financial gain. However, the existence of prosocial preferences has been inferred post hoc from the results of economic games, rather than with direct experimental tests. Here, we test how behavior in a public-goods game is influenced by knowledge of the consequences of actions for other players. We found that (i) individuals cooperate at similar levels, even when they are not informed that their behavior benefits others; (ii) an increased awareness of how cooperation benefits others leads to a reduction, rather than an increase, in the level of cooperation; and (iii) cooperation can be either lower or higher than expected, depending on experimental design. Overall, these results contradict the suggested role of the prosocial preferences hypothesis and show how the complexity of human behavior can lead to misleading conclusions from controlled laboratory experiments.

altruism | behavioral economics | black box | framing effect | reciprocity

Economic experiments have shown that people cooperate at levels higher than predicted if they were maximizing their financial gain (1–6). For example, in public-goods games, which are used as a model for collective-action problems, individuals can contribute money to produce benefits that are shared by all members of their group, including themselves. If an individual's share of the public benefit from his contribution is less than his contribution, then any individual wishing to maximize his own income should contribute zero. In contrast to this prediction, most individuals do contribute something (typically around 40–50%), and although contributions go down when play is repeated, groups continue to contribute 10–20% of their resources to the public good (4, 7).

These results have led to the argument that human behavior takes into account the welfare of others in a way that is perhaps unique in the animal world and cannot be explained by standard gene-based evolutionary theory (2–5, 8). Specifically, it is argued that individuals have prosocial preferences that lead them to help others, even when interacting with unrelated strangers, and when there is no scope of repeated interactions (2–5, 8). However, the existence of prosocial preferences has only been inferred post hoc from the existing data, rather than tested for directly, by varying among treatments the extent to which individuals can discern the consequences of their behavior for others (9–15).

Here, we directly test the prosocial preferences hypothesis by experimentally varying the information that individuals receive about the consequences of their behavior for others. We used a linear public-goods game, with groups of four players, in which contributions to the public good were multiplied by 1.6 before being shared out equally. This means that for each monetary unit (MU) contributed a player received 0.4 units back (the marginal per-capita return, MPCR), and so the strategy that maximized the financial gain of individuals was to contribute 0 MU (0% cooperation) from their endowments of 40 MU (3). We randomly determined group membership for each of 20 rounds to reduce the possibility for reciprocity while allowing for learning (16).

We carried out three versions of this game. The only difference between versions was the information that individuals received about the consequences of their behavior. This information was either standard, reduced, or increased (Table 1 and Fig. S1). In our standard-information treatment we told participants the payoff structure of the game and then, after each round of play, how much the others in their group had contributed and their own personal payoffs. This is the same level of information provided in many previous studies (3). To provide reduced information, we used a “black-box” treatment, in which participants were not even informed that they were playing with others in a group. Specifically, we told participants that they could input 0–40 virtual “coins” into a virtual black box, which would return a (nonnegative) output, determined by a mathematical function on their input. To provide an “enhanced-information” treatment, we carried out the same public-goods game but also provided individuals with a detailed breakdown of how much each other member of their group had contributed, received back from the public good, and earned in that round (Table 1 and Fig. S1).

These three treatments allowed us to test whether prosocial preferences are either necessary or sufficient to explain the higher-than-expected level of cooperation in public-goods games. In our black-box treatment, players could only learn the payoff-maximizing strategy (contribute 0 MU) through trial and error, and they did not even know that they were playing a game with other individuals, let alone that their contributions had beneficial consequences for others. Consequently, prosocial preferences cannot be invoked to explain the data. In contrast, in our enhanced-information treatment, players were still given the same information about the game payoffs, but they were better able to see that their cooperation was costly to themselves and beneficial to others. The prosocial preferences hypothesis predicts that this should either not alter cooperation levels or lead to relatively higher levels of cooperation. Furthermore, as an experimental control, we repeated all three of our above treatments but with a higher benefit from contributing to the public good, such that the strategy that would maximize financial gain was to cooperate 100% (contribute 40 MU).

Results

In total, we had 16 sessions and used 236 participants, who all played both the black-box game and then, or previously, also participated in either the standard-information or the enhanced-information public-goods game, which they were informed was a separate experiment. Each treatment was played for 20 successive rounds. Overall, there was a significant difference in the mean

Author contributions: M.N.B.-C. and S.A.W. designed research, performed research, analyzed data, and wrote the paper.

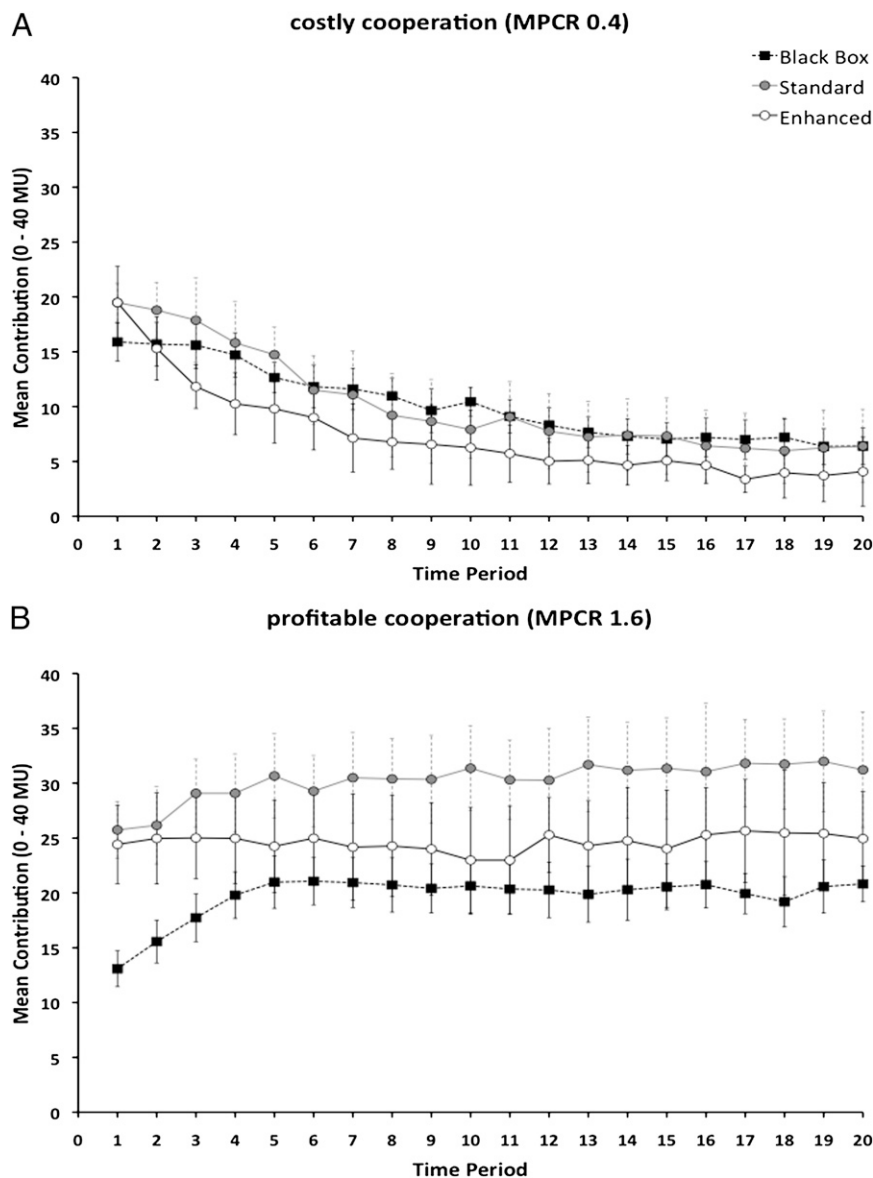
The authors declare no conflict of interest.

*This Direct Submission article had a prearranged editor.

Freely available online through the PNAS open access option.

¹To whom correspondence should be addressed. E-mail: stuart.west@zoo.ox.ac.uk.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1210960110/-DCSupplemental.



cooperation. The prosocial preferences hypothesis would suggest that increased information would either not influence or would increase the level of cooperation. If players' preferences are being measured by the consequences for others resulting from their decisions, then this assumes that they are aware of the effects of their decisions. Therefore our increased information should have no effect on contribution levels as it is only telling players what they already know (2). Alternatively, if the players are perhaps incapable of calculating the consequences of their behavior for others, but do cooperate because of prosocial preferences, then this information would lead to increased cooperation levels. Consequently, not only do our results fail to provide experimental support for prosocial preferences, they suggest the opposite: participants cooperate less when they have clearer information about how their cooperation benefits others. Instead, the observed pattern again supports the null hypothesis, because the increased information always suggests that contributions are costly and thus may influence uncertain players to play correctly in the costly game and play incorrectly in the profitable game (from an income-maximizing perspective).

Third, prosocial preferences cannot explain our result that, in our profitable public-goods game treatments, the mean contributions in all rounds remained substantially below 40 MU (100% cooperation; Fig. 1B). The prosocial preferences hypothesis predicts 100% cooperation (40 MU) in this profitable game. Alternatively, it might be argued that rational individuals prioritizing within-group success could predict <100% cooperation in profitable games (Fig. 1B). However this would require participants to value their relative success within a temporary, transient group (we changed our group compositions randomly each round) more than their relative success within the session or the experiment as a whole (those who contribute more make more in the session). The general point here is that because previous experiments lacked appropriate controls for imperfect behavior, any imperfect behavior would have led to a higher-than-expected level of cooperation, and hence a biased conclusion.

Alternative Explanations. Could our results have been influenced by the potential for reciprocity (16)? Our participants played for 20 rounds and thus could expect some repeated interactions. However, this is likely to have had negligible influence for three reasons. First, group formation was random and anonymous, such that there was no potential for reputational consequences and one could not predict when a repeated interaction would occur. Second, there was the same amount of repeated interaction in all treatments, and hence this could not explain the between-treatment differences upon which we focus. Third, our data follow the normal, stylized patterns of public-goods games, either with or without repeated interactions, and it has been shown elsewhere that with random group composition there was no significant influence of whether there could be some repeated interactions (21).

Could it be argued that the black-box treatment is sufficiently similar to a social interaction to trigger any purported human tendencies to altruistically help others, and that this explains their pattern of monetary contributions? Such an explanation could explain the previous result that people "cooperate" when knowingly playing public-goods games with a computer (22). However, our black-box game was deliberately devoid of all social cues and frames, and so it is harder to envisage how this could have engaged social tendencies. Furthermore, if we accept that people cooperate maladaptively with computers because of hard-wired tendencies to be prosocial, then we would have to apply the same logic to all cooperation in economic games. This would mean that the cooperation observed in public-goods games would be explained by hard-wired responses to evolutionary factors such as the consequences of repeated interactions and reputational concerns, even when these are removed by the experimenter, and thus the prosocial preferences explanation for costly cooperation with strangers is invalidated either way (15, 16).

Could it be argued that some form of social preference, such as reciprocity or inequality aversion, can explain both the decline in cooperation over time in costly public-goods games and the lower cooperation in the enhanced-information treatment (Fig. 1A)? The idea here is that both of these mechanisms would place a negative value on increased earnings by others, and increased information would lead to individuals' cooperating less to either punish nonreciprocators or reduce the inequality (2, 23, 24). However, these hypotheses are clearly falsified by the fact that in our profitable control treatment the mean level of cooperation increased over time, and from an intermediate start that is not predicted by such preferences. Furthermore: (i) as described above, our experiment was designed to minimize the potential for reciprocity; (ii) as discussed in greater detail below, the inequality aversion hypothesis often implicitly assumes that individuals compute the financial consequences of their strategy from the game structure—in this case, the enhanced-information would provide no new insight, and so there should be no influence of enhanced-information (2); and (iii) our black-box treatment shows that the assumption that the decay in public-goods games is due to "punishing" others or reducing inequality is not required, because this decline still occurs when participants can only be responding to their payoffs (Fig. 1A).

Implications. Given that the prosocial preferences hypothesis appears neither a necessary nor sufficient explanation, how can we explain the higher-than-expected level of cooperation that is often observed in economic games? The most parsimonious explanation for our results is our null hypothesis presented above, that (i) individuals are trying to maximize their financial gain, but (ii) behavior is imperfect due to uncertainty or false beliefs, or subject to some sort of noise, which could result from a variety of factors, including errors, boredom, learning, exploration, fluctuating preferences, or evolutionary constraints. This combination of factors can explain three results that are inconsistent with prosocial preferences: (i) higher-than-expected levels of cooperation were seen when the income-maximizing strategy was to contribute zero (0%; Fig. 1A), but lower-than-expected levels of cooperation were seen when the income-maximizing strategy was to cooperate completely (100%; Fig. 1B) (17); (ii) the level of cooperation in a standard-information public-goods game did not differ significantly from that in a black-box game in which individuals had no information about the consequences of their behavior for others (Fig. 1A); and (iii) enhanced information about the earnings of others, which suggests that contributions are costly, led to lower, not higher, levels of cooperation regardless of whether cooperation was costly or profitable (Fig. 1). Although we have not directly tested the role of punishment, our results suggest the possibility of analogous problems for work on punishment, because this previous work has also assumed that the higher-than-expected levels of cooperation were because of prosocial preferences (3).

More generally, our results show how the complexity of human behavior can lead to potentially misleading conclusions from controlled laboratory settings. Even though there is evidence to suggest that many people fail to fully understand public-goods games (17, 22, 25), rational choice theory has been extended to include prosocial preferences in a post hoc attempt to explain the data (2, 8, 26): "Expanding the domain of preferences to include the utility of others provides a coherent way to extend rational choice theory" (ref. 26, abstract). Thus, the paradigm of prosocial preferences has been based on the assumption that people act rationally (perfectly) in accord with their desires and the financial consequences of their decisions (2, 8, 26). In the extreme, it has been argued that "it is not necessary for subjects to be informed about the final monetary payoffs of other subjects" and that all that is required for relative payoffs to determine behavior is knowledge of the game structure such that individuals can

“compute the distributional implications” (ref. 2, footnote 6). However, our results show that humans do not behave in such a “robotlike” way and thus may not perfectly express their preferences, especially when these preferences are “measured” by the social consequences (relative payoffs) of their actions in one-shot experimental games (6, 13, 15, 27, 28). Furthermore, the level of cooperation observed is “irrationally” sensitive to parameters that do not change the rational strategy to fully defect, such as group size, group comparisons, and the presence of eyelike images (29–33). Whereas such behavior can be beneficial in the real world (15, 16, 34), it can lead to seemingly irrational behavior in experimental settings, where normally useful cues of success can be misleading, and the influence of usually important factors has been experimentally excluded. This does not mean that humans lack a sense of fairness, but rather that it is most favored when it provides a greater net benefit.

Methods

Participants and Sessions. A total of 236 voluntary participants (108 females, 111 males, 17 unknown) took part in one session each. Each session lasted ~50–60 min. We had 12 or 16 participants for each session (5 and 11 sessions, respectively) and conducted all 16 sessions at the Centre for Experimental Social Sciences (CESS) at Nuffield College, University of Oxford. CESS recruited the participants using the Online Recruitment System for Economic Experiments (35) and upon our request only recruited participants that had never before participated in a public-goods experiment (at CESS). Although we did not fully inform participants about the consequences of their decisions in the black-box treatments, we did not lie to them, and the CESS ethical committee, which forbids deception in economic experiments, approved the experiment. We programmed and conducted the experiment in z-Tree (36). The decisions and responses of the participants were anonymous, as were their earnings and payments, which were administered by the CESS administrators, who were not directly involved in the experiments. In total, we gave each participant a sum total endowment of £4.80 to allocate to either her private account or her group account/black box. The mean monetary reward was £12.40 and ranged from £6.20 to £15.50.

Different Treatments. An overview of the information we provided as feedback during the treatments is shown in Fig. S1. From the participants' points of view, they played two separate experiments (the black-box experiment and the public-goods experiment), each of them under two conditions for 20 times apiece. The two conditions were (i) cooperation is directly costly and (ii) cooperation is directly profitable, although participants were not necessarily aware of this distinction. The order of the treatments was counterbalanced across sessions, but the two black-box games (costly and profitable) were always the first two or the final two games of the session. Within this constraint, we used all possible treatment order permutations ($n = 2 \times 2 \times 2 = 8$) to produce an orthogonal experimental design (Table S1).

Black-box experiment. All participants played the same two versions of the black-box game. In short, we told the participants that they had received 40 virtual coins, worth real money, and that they had to decide to either keep all 40 coins or to input a fraction (0–40 coins) of them into a “black box.” They then read the following information (a full copy of the instructions is given in SI Appendix): “This ‘black box’ performs a mathematical function that converts the number of coins inputted into a number of coins to be outputted. The function contains a random component, so if two people were to put the same amount of coins into the ‘black box’ they would not necessarily get the same output. The number outputted may be more or less than the number you put in, but it will never be a negative number, so the lowest outcome possible is to get 0 (zero) back. If you choose to input 0 (zero) coins, you may still get some back from the box.”

We also told them that they would play for 20 turns, receiving a new set of 40 coins each time, and that the most they could ever input was 40 coins. We also explained that after their 20 turns they would start again with a new, potentially different black box for another 20 turns. In addition, we explained that their final income from each turn would be their initial 40 coins, minus their input, plus all of the coins, if any, that they get back. The feedback we gave them after each turn comprised four items: “initial coins” (i.e., their endowment), “input” (i.e., their contribution), “output returned” (i.e., their returns from the group project), and their final balance from that round (Fig. S1).

Statistical Analysis. Because the responses of participants within a session are nonindependent we used the mean level of contributions for each session per time period to avoid pseudoreplication. We used LMM and fitted “session ID” as a random factor to account for the repeated measures over time. When comparing the mean levels of contributions for the whole treatment, we excluded the first time period of each treatment. This was because the treatment effect was not applied until after the first contributions had been made, and thus any differences in the first time period could only be due to chance. We compared the relative merit of different models by using maximum likelihood models and comparisons of the residual deviance. The comparison is made with a χ^2 difference test of the change in deviance between models, with the degrees of freedom set by the difference in the number of parameters between two nested models. Models are progressed from the null model through to the maximal model, and vice versa (although the order is unimportant with orthogonal experimental designs such as ours), in search of the minimal adequate model (MAM). The MAM is the model with the smallest residual deviance for a given level of parameters and is identified on the basis of deletion/addition tests. These are χ^2 tests that assess the significance of the change in deviance that results when a given term is removed/inserted. The MAM is the model that can no longer be improved by the insertion of additional terms or factor levels (37).

ACKNOWLEDGMENTS. We thank Luis Miller at the Centre for Experimental Social Sciences for recruiting participants, the reviewers for their helpful comments, Nuffield College for facilities, and the European Research Council for funding.

1. Camerer CF, Fehr E (2006) When does “economic man” dominate social behavior? *Science* 311(5757):47–52.
2. Fehr E, Schmidt KM (1999) A theory of fairness, competition, and cooperation. *Q J Econ* 114(3):817–868.
3. Fehr E, Gächter S (2002) Altruistic punishment in humans. *Nature* 415(6868):137–140.
4. Fehr E, Fischbacher U (2003) The nature of human altruism. *Nature* 425(6960):785–791.
5. Gintis H, Bowles S, Boyd R, Fehr E (2003) Explaining altruistic behavior in humans. *Evol Hum Behav* 24(3):153–172.
6. Henrich J, et al. (2005) “Economic man” in cross-cultural perspective: Behavioral experiments in 15 small-scale societies. *Behav Brain Sci* 28(6):795–815, discussion 815–755.
7. Ledyard J (1995) Public goods: A survey of experimental research. *Handbook of Experimental Economics*, eds Kagel J, Roth A (Princeton Univ Press, Princeton), pp 253–279.
8. Fehr E, Fischbacher U (2005) Human altruism: Proximate patterns and evolutionary origins. *Anal Kritik* 27:6–47.
9. Binmore K, Shaked A (2010) Experimental economics: Where next? *J Econ Behav Organ* 73(1):87–100.
10. Blanco M, Engelmann D, Normann HT (2011) A within-subject analysis of other-regarding preferences. *Games Econ Behav* 72(2):321–338.
11. Burnham TC, Johnson DDP (2005) The biological and evolutionary logic of human cooperation. *Anal Kritik* 27:113–135.
12. Harrison GW (2008) Neuroeconomics: A critical reconsideration. *Econ Philos* 24(3):303–344.
13. Oechssler J (2003) Intentions matter: Lessons from bargaining experiments — Comment on, reasons for conflict: Lessons from bargaining experiments, by Falk, A., Fehr, E., and Fischbacher, U. *J Inst Theor Econ* 159(1):195–198.
14. Ross D (2012) What can economics contribute to the study of human evolution? *Biol Philos* 27:287–297.
15. Trivers RL (2004) Mutual benefits at all levels of life. *Science* 304(5673):964–965.
16. Trivers RL (1971) Evolution of reciprocal altruism. *Q Rev Biol* 46(1):35.
17. Kümmerli R, Burton-Chellew MN, Ross-Gillespie A, West SA (2010) Resistance to extreme strategies, rather than prosocial preferences, can explain human cooperation in public goods games. *Proc Natl Acad Sci USA* 107(22):10125–10130.
18. Hagen EH, Hammerstein P (2006) Game theory and human evolution: A critique of some recent interpretations of experimental games. *Theor Popul Biol* 69(3):339–348.
19. Engelmann D, Strobel M (2010) Inequality aversion and reciprocity in moonlighting games. *Games* 1:459–477.
20. Young HP (2009) Learning by trial and error. *Games Econ Behav* 65(2):626–643.
21. Fehr E, Gächter S (2000) Cooperation and punishment in public goods experiments. *Am Econ Rev* 90(4):980–994.
22. Houser D, Kurzban R (2002) Revisiting kindness and confusion in public goods experiments. *Am Econ Rev* 92(4):1062–1069.
23. Levine DK (1998) Modeling altruism and spitefulness in experiments. *Rev Econ Dyn* 1(3):593–622.
24. Fehr E, Gächter S (2000) Fairness and retaliation: The economics of reciprocity. *J Econ Perspect* 14(3):159–181.
25. Andreoni J (1995) Cooperation in public-goods experiments: Kindness or confusion? *Am Econ Rev* 85(4):891–904.

26. Sobel J (2005) Interdependent preferences and reciprocity. *J Econ Lit* 43(2):392–436.
27. Haselton MG, Nettle D (2006) The paranoid optimist: An integrative evolutionary model of cognitive biases. *Pers Soc Psychol Rev* 10(1):47–66.
28. McCullough ME, Kimeldorf MB, Cohen AD (2008) An adaptation for altruism? The social causes, social effects, and social evolution of gratitude. *Curr Dir Psychol Sci* 17(4):281–285.
29. Isaac RM, Walker JM (1988) Group-size effects in public-goods provision: The voluntary contributions mechanism. *Q J Econ* 103(1):179–199.
30. Burton-Chellew MN, West SA (2012) Pseudocompetition among groups increases human cooperation in a public-goods game. *Anim Behav* 84:947–952.
31. Haley KJ, Fessler DMT (2005) Nobody's watching? Subtle cues affect generosity in an anonymous economic game. *Evol Hum Behav* 26(3):245–256.
32. Bateson M, Nettle D, Roberts G (2006) Cues of being watched enhance cooperation in a real-world setting. *Biol Lett* 2(3):412–414.
33. Burnham TC, Hare B (2007) Engineering human cooperation: Does involuntary neural activation increase public goods contributions? *Hum Nature-Int Bios* 18(2):88–108.
34. Trivers RL (2006) Reciprocal altruism: 30 years later. *Cooperation in Primates and Humans: Mechanisms and Evolution*, eds Kappeler PM, van Schaik CP (Springer, Berlin).
35. Greiner B (2004) The online recruitment system ORSEE 2.0: A guide for the organization of experiments in economics. *Working Paper Series in Economics* 10 (Department of Economics, University of Cologne, Cologne, Germany).
36. Fischbacher U (2007) z-Tree: Zurich toolbox for ready-made economic experiments. *Exp Econ* 10(2):171–178.
37. Crawley MJ (2007) *The R Book* (Wiley, Chichester, UK), p viii.