

# Habits of Virtue: Creating Norms of Cooperation and Defection in the Laboratory<sup>1</sup>

Alexander Peysakhovich<sup>2</sup> & David G. Rand<sup>3</sup>

*Forthcoming in Management Science*

What explains variability in norms of cooperation across organizations and cultures? One answer comes from the internalization of norms prescribing behavior that is typically successful under the institutions that govern one's daily life. These norms are then carried over into atypical situations beyond the reach of institutions. Here we experimentally demonstrate such spillovers. First, we immerse subjects in environments that do or do not support cooperation using repeated Prisoner's Dilemmas. Afterwards, we measure their intrinsic prosociality in one-shot games. Subjects from environments that support cooperation are more prosocial, more likely to punish selfishness, and more generally trusting. Furthermore, these effects are most pronounced among subjects who use heuristics, suggesting that intuitive processes play a key role in the spillovers we observe. Our findings help to explain variation in one-shot anonymous cooperation, linking this intrinsically motivated prosociality to the externally imposed institutional rules experienced in other settings.

---

<sup>1</sup> We thank Robert Aumann, Colin Camerer, John Clithero, Armin Falk, Guillaume Frechette, Drew Fudenberg, Ed Glaeser, Joseph Henrich, Benedikt Herrmann, Moshe Hoffman, Jillian Jordan, David Laibson, Martin Nowak, Nathan Nunn, Pietro Ortoleva, Aurelie Ouss, Antonio Rangel, Peter Richerson, Daria Roithmayer, Al Roth, Klaus Schmidt, Dmitry Taubinsky, Julian Wills, members of the Rangel Lab and seminar participants at Harvard, Princeton, Yale, MIT, Brown and Stony Brook for their invaluable comments, and the John Templeton Foundation for providing funding for this work. The authors declare no conflict of interest.

<sup>2</sup> [alex.peys@gmail.com](mailto:alex.peys@gmail.com); Program for Evolutionary Dynamics, Harvard University, Cambridge, MA, and Department of Psychology, Yale University, New Haven, CT

<sup>3</sup> [david.rand@yale.edu](mailto:david.rand@yale.edu); Department of Psychology, Department of Economics, School of Management, Yale University, New Haven, CT

## 1. Introduction

The tendency to cooperate, or to pay a personal cost to give a benefit to another person or group, is not a constant between groups: norms of cooperation and trust vary markedly across organizations (Leana and Buren, 1999, McAllister, 1995, Rousseau et al., 1998), countries (Cappelen et al., 2013, Ellingsen et al., 2012, Henrich et al., 2010, Herrmann et al., 2008) and cultures (Gächter et al., 2010, Henrich et al., 2005, Sapienza et al., 2006). Here we use economic game experiments to shed light on the origins of this heterogeneity. We argue that cooperative norms (i.e. individual-level conceptions of appropriate behavior)<sup>4</sup> and expectations about the behavior of others are driven, at least in part, by spillovers from representative daily life interactions.<sup>5</sup> Individuals who primarily interact in environments where the “rules of the game” make cooperation advantageous get in the habit of cooperating, and as a result are more cooperative even in one-shot interactions (such as those in most laboratory games). By this logic, institutional differences across groups lead to differences in cooperative norms between those groups<sup>6</sup>: members of groups with institutions that successfully incentivize cooperation most of the time continue to behave more cooperatively even in interactions where no such institutional incentives exist.

In this paper we provide empirical evidence for this line of reasoning with a set of experiments directly demonstrating how randomly assigned institutions can change norms and create “cultural differences” in the lab. Each experiment has the same basic structure: first, subjects are assigned to interact in an environment where, because of the rules of interaction, cooperative equilibria either are strongly supported or do not exist (Stage A). After experience in this environment, subjects take part in a battery of one-shot anonymous economic games used to study cooperation (Stage B). We then look for the impact of random assignment to experimental environment in Stage A on behavior in Stage B.

---

<sup>4</sup> There are many different uses of the word “norm” in different literatures. Here, we use “norm” to refer to an individual’s internalized conception of what behavior is appropriate in a given setting or range of settings (manifested in terms of that individual’s preferences, in addition to beliefs). By this definition, holding a particular norm makes one feel personally compelled to engage in behaviors that the norm deems to be appropriate, and also makes one upset when other people do not follow the norm and instead engage in behaviors seen as inappropriate (Fehr and Fischbacher, 2004a, Fehr and Fischbacher, 2004b, Jordan et al., 2014). For example, individuals with a norm prescribing cooperation will be both inclined to cooperate and to punish non-cooperators, even at a cost. Thus these norms form the underpinnings of social preferences.

<sup>5</sup> Spillover effects occur when subjects extrapolate from experience in one domain to guide behavior in other domains with different incentive structures (e.g. Rand et al. (2014c)). This includes generalization that occurs outside the lab, where strategies from typical settings with future consequences are applied to atypical situations where risk-free exploitation is possible, as well as generalization that occurs inside the lab, either from typical settings outside the lab to one-shot anonymous lab games, or from repeated games in the lab to subsequent one-shot games in the lab (as in our experiments).

<sup>6</sup> We use a broad definition of “institution” that encompasses any factor external to the individual that affects incentives (in contrast to norms, which are internalized conceptions). By this definition, repeated interactions, reputation systems, and the threat of punishment by third parties or organizations are all forms of institutions.

To create Stage A lab environments that favor cooperation or non-cooperation, we use infinitely repeated Prisoner's Dilemma games. The extent to which an environment supports cooperation can be manipulated using many mechanisms other than repeated interactions, including reputation, sanctions, partner choice, intergroup competition, and formal institutions (for a review, see (Jordan et al., In press)). We chose repeated games as a model of the future consequences created by these mechanisms because the determinants of the emergence (and stability) of cooperation in repeated games are well understood both experimentally (Dal Bó, 2005, Dal Bó and Fréchet, 2011, Fudenberg et al., 2012, Rand and Nowak, 2013) and theoretically (Blonski et al., 2011, Fudenberg and Maskin, 1986, Fudenberg and Maskin, 1990, Mailath and Samuelson, 2006). This allows us to select combinations of payoffs and continuation probabilities that lead to high levels of cooperation in one treatment (the 'C-Culture' treatment) and low levels of cooperation in the other (the 'D-Culture' treatment).

In Experiment 1, the repeated Prisoner's Dilemmas of Stage A are followed by a Stage B consisting of one-shot anonymous cooperation games: the Public Goods Game, Trust Game, Dictator Game and Ultimatum Game. Each of these games involves one or more choices of whether to transfer money from oneself to one or more others (i.e. to act prosocially). We show that subjects randomized into the C-Culture in Stage A are substantially more prosocial in the Stage B games compared to subjects randomized into the D-Culture treatment. Importantly, this is true even in the Dictator Game, where the recipient is passive and therefore expectations about the decisions of others play no role. We also provide several pieces of evidence that expectations about the *type* (i.e. level of altruism) of one's Stage B co-player do not explain our treatment effect. One such piece of evidence comes from replicating our treatment's effect on the DG in a Supplemental Experiment using an online non-student subject pool where Stage B recipients are complete strangers, do not participate in Stage A, and take no actions in the experiment whatsoever (other than receiving any money given to them). This replication shows that the treatment effect cannot be explained by reciprocity towards the individuals one interacted with in Stage A, and speaks against inferring that one's DG recipient would have been more or less altruistic based on the observed level of cooperation in one's Stage A partners. Thus, in these experiments, we demonstrate that the "rules of the game" in one setting can strongly influence behavior in other settings with different incentives.

We also explore the cognitive mechanism through which these spillovers occur. To do so, we take a dual-process perspective and conceptualize decision-making as the result of interactions between intuitive and deliberative processes (Epstein et al., 1996, Gilovich et al., 2002, Kahneman, 2003, Tversky and Kahneman, 1974). Intuitive processes are fast, automatic, emotional and heuristic in nature: intuitions often favor behaviors which are advantageous in typical settings (Gigerenzer and Goldstein, 1996, Gigerenzer et al., 1999). Deliberation, by contrast, is slow, controlled, rational, and tailored to the specific

decision under consideration. Thus in *atypical* settings, deliberation can override sub-optimal intuitive, heuristic responses (which, although they may be optimal in general, are ill-matched to the atypical setting at hand).

In previous work, this dual process lens has been applied to cooperation via the Social Heuristics Hypothesis (SHH) (Rand, et al., 2014c), which takes theories of cultural evolution and norm internalization (Bowles and Gintis, 2002, Bowles and Gintis, 2003, Boyd and Richerson, 2009, Chudek and Henrich, 2011, Gintis, 2003, Henrich et al., 2006, Richerson and Boyd, 2005) and makes them explicitly dual process. The SHH contends that the internalization of norms occurs via the channel of intuitive processing. By this account, social behaviors that are rewarded in the course of one's daily life (in the developed world, this is typically cooperation) become internalized as default heuristics. In one-shot anonymous interactions (for example, in the lab), deliberation then leads one to realize that selfishness is actually optimal. Empirical evidence for this theory comes from the finding that experimentally inducing heuristic processing (via time pressure, conceptual priming, or cognitive load) can increase prosociality in economic games (Cone and Rand, 2014, Cornelissen et al., 2011, Lotz, 2014, Rand et al., 2012, Rand et al., 2014b, Rand, et al., 2014c, Roch et al., 2000, Schulz et al., 2014) (although other studies have found null effects, e.g. Hauge et al. (2014), Tinghög et al. (2013), Verkoeijen and Bouwmeester (2014)).

In the context of the present experiments, therefore, the SHH predicts that experiences in Stage A affect behavior in Stage B by remodeling subjects' default responses; and therefore that Stage A should have less of an effect on Stage B behavior among subjects that are better at overriding their intuitive heuristic-based responses. Consistent with this prediction, we find that subjects in Experiment 1 who engage in more deliberative thinking (as measured by the 'Cognitive reflection test' (Frederick, 2005)) are much less influenced by Stage A when making their Stage B decisions.

In Experiment 2, we ask whether exposure to a cooperative environment in Stage A also affects the willingness to *enforce* cooperation. To investigate this issue, we follow Stage A with a Stage B that consists of punishment games where subjects can pay to reduce the earnings of others (Fehr and Fischbacher, 2004b, Fehr and Gächter, 2000, Ostrom et al., 1992, Ouss and Peysakhovich, 2013). We find that subjects randomized into the C-Culture in Stage A are more likely to engage in third party punishment of selfishness than those randomized into the D-Culture; and find some evidence that subjects in the D-Culture may be more likely to engage in anti-social punishment of cooperators in a Public Goods Game. We also replicate the result from Experiment 1 that, as predicted by the SHH, Stage A has a much smaller effect on Stage B behavior among subjects who engage in more deliberative thinking. Thus, in Experiment 2, we provide evidence that Stage A affects *norms* of cooperation, altering enforcement behavior as well as cooperative choice.

Finally, we ask whether our Stage A manipulation has effects that generalize beyond behavior in economic games. To do so, we have subjects answer a commonly used question regarding generalized trust from the World Values Survey, which has been shown to reflect prosocial orientation (not just beliefs about the trustworthiness of others). Consistent with our Stage B game results, we find that subjects randomized into the D-Culture treatment in Stage A report being significantly less generally trusting of others than those from the C-Culture, and that this effect is more pronounced among subjects who rely on heuristics.

Thus, by exposing our subjects to environments that favor cooperation or defection, we replicate three major results in cross-cultural studies within a single subject pool: variation in prosociality (Henrich, et al., 2005, Henrich, et al., 2010), variation in norm-enforcement (Ellingsen, et al., 2012, Gächter, et al., 2010, Henrich, et al., 2005, Henrich, et al., 2010, Henrich, et al., 2006, Herrmann, et al., 2008) and variation in generalized trust (La Porta et al., 2001, Putnam, 2000). In doing so, we provide causal evidence in support of previous cross-cultural correlational studies suggesting that norms related to cooperation (revealed by play in one-shot economic games) positively co-vary across cultures with measures of institutional quality such as rule of law (Gächter, et al., 2010, Herrmann, et al., 2008) and market integration (Henrich, et al., 2010)<sup>7</sup>: we show that a qualitatively equivalent pattern of results can be generated through random assignment in the laboratory, operating through the channel of heuristic decision-making. Thus our experiments support the argument that institutional differences across organizations and cultures can affect internalized norms of prosociality, and demonstrate how norms can be changed through a top-down process driven by institution designers.

Our paper proceeds as follows. In Section 2, we describe the design of Stage A, which is the same in both Experiment 1 and Experiment 2. In Section 3, we describe the design of Experiment 1's Stage B, which investigates prosociality, and present the results of Experiment 1. In Section 4, we describe the design of Experiment 2's Stage B, which investigates punishment behavior, and present the results of Experiment 2. In Section 5, we aggregate across studies and investigate treatment effects on generalized trust (rather than game play). In Section 6, we present a concluding discussion.

## **2. General experimental design for creating cultures of cooperation or defection (Stage A)**

Our two-stage experiments are designed to evaluate the effect of exposure to environments that incentivize cooperation or non-cooperation in Stage A on subsequent behavior in one-shot anonymous

---

<sup>7</sup> (Bowles, 1998) surveys more evidence to this effect and presents a model in which institutional factors can affect the endogenous evolution of preferences.

interactions in Stage B. In Stage A, subjects play a series of stochastically repeated Prisoner’s Dilemma (RPD) games with different partners. At the beginning of each game, subjects are matched in pairs. Each game consists of a random number of rounds. In each round, subjects play a simultaneous PD stage game: both players choose an action, C or D (labeled A and B for the subjects). Subjects are then informed of the decision of their partner, and the resulting earnings of each player for the round. They then play another round with the same partner with probability  $\delta$ , while with probability  $1-\delta$  the game ends and subjects are rematched with a new partner for a new game (and informed of this rematching).

To manipulate the extent to which the Stage A environment supports cooperation versus non-cooperation, subjects are randomly assigned to one of two treatments: the C-Culture treatment or the D-Culture treatment. The treatments differ both in continuation probabilities, with  $\delta = 7/8$  (expected game length of 8 rounds) in the C-Culture and  $\delta = 1/8$  (expected game length of 1.14 rounds) in the D-Culture, and the PD game payoffs, with the D-Culture involving a higher temptation payoff from defection exploiting cooperation. The stage game payoffs (shown in Monetary Units, MU) were

<i>C-Culture</i>	C	D	<i>D-Culture</i>	C	D
C	4, 4	0, 5	C	4, 4	0, 6
D	5, 0	1, 1	D	6, 0	1, 1

At the end of the experiment, MU are converted to cash at an exchange rate of \$1 = 30 MU.

An important determinant of whether cooperation emerges in infinitely repeated games is whether the strategy Tit-for-Tat (TfT) risk dominates Always Defect (AllD) (Blonski, et al., 2011, Rand and Nowak, 2013). For our C-Culture we therefore choose parameters such that TfT strongly risk dominates AllD, and so we expect subjects in the C-culture to learn to cooperate and maintain cooperation in the long-run.<sup>8</sup> In the D-Culture treatment, conversely, we choose a specification such that cooperation is not an equilibrium, and therefore we expect subjects to learn to defect<sup>9</sup>. Thus Stage A should create environments in which players consistently cooperate (C-Culture) or consistently defect (D-Culture).

To control for random variation in lengths between sessions, we follow the procedure of Dreber et al. (2008), Fudenberg, et al. (2012), Rand et al. (2014a). For each of the two treatments, we generate a single set of game lengths using the appropriate distribution, and then use this same set of game lengths in

<sup>8</sup> The payoff of TfT against a 50-50 mix between Tft and AllD is given by  $.5 (4)*(1/(1-\delta)) + .5 (0 + 1*\delta/(1-\delta)) = 19.5$ , whereas the payoff of AllD against the same mix is given by  $.5 (5) + .5 (1) + 1*(\delta/(1-\delta)) = 10$ . Thus TfT strongly risk-dominates AllD.

<sup>9</sup> No cooperative equilibria exist as even in the presence of the harshest possible punishment (Grim Trigger), the present gain from defecting (2 for sure) outweighs the expected future losses from loss of cooperation ( $3 * (\delta/(1-\delta)) \sim .33$ ).

every session. The specific game lengths used are shown in the Appendix Tables A18 and A19. In Experiment 1, subjects play a total of 53 rounds (split into 10 games) in the C-Culture treatment, and 51 rounds in D-Culture treatment (split into 45 games). In Experiment 2, the Stage B decisions take longer to complete, and so we shorten the total length of Stage A to be approximately 40 PD rounds in total (split into 7 games in the C-Culture and 35 in the D-Culture).

Our goal is to examine the consequences of the resulting habituation and reinforcement of either cooperation or defection on subsequent behavior. A potential confound, however, exists in the form of income effects: if subjects in the C-Culture treatment cooperate much more than those in the D-Culture treatment, Stage A earnings will be higher in the C-Culture treatment. To ensure that any differences we observe in Stage B result from acculturation to cooperation or defection rather than any consequences of differences in income, we vary the size of the initial endowment subjects receive at the beginning of Stage A: in the D-Culture treatment, subjects begin with an endowment of 150 MU, whereas in the C-Culture treatment subjects begin with only 40 MU (50 MU in Experiment 2 to control for shorter Stage A length). As a result, subjects in each condition finish Stage A with similar earnings on average.<sup>10</sup> Thus, our results cannot be explained by players in the C-Culture earning more than players in the D-Culture.

After completing Stage A, subjects enter Stage B and play a battery of one-shot anonymous games commonly used to assess prosociality (Experiment 1) or punishment (Experiment 2). Subjects are given no information about the history of Stage A play of their Stage B interaction partners. The Stage B games are described in more detail in the corresponding sections for each experiment below.

Unless otherwise noted, all analyses reported use linear regression. To correct for within-session correlation induced by Stage A behavior, we cluster standard errors at the session level. To address potential small sample bias due to the low number of sessions, we calculate significance levels for each regression by bootstrapping, with 1,000 iterations per bootstrap. This clustering also addresses correlation across decisions made by a particular subject, and replacing session-level clustering with subject-level clustering calculated via standard normal approximations does not qualitatively change any of our results.

Our experiments were run in the Harvard Decision Sciences Laboratory between April and September 2012, and implemented using the Z-tree software (Fischbacher, 2007). We aimed for approximately 50 subjects per condition, as per the recommendations of Simmons et al. (2013).

---

<sup>10</sup> Subjects in the C-Culture earned substantially more MU during the RPD than subjects in the D-Culture (Experiment 1: C-Culture 166.18 MU vs D-Culture 80.65 MU; Experiment 2: C-Culture 110.25 MU vs D-Culture 61.68 MU). However, the difference in initial endowments more than made up for this difference. Thus, upon leaving Stage A, subjects in the C-Culture had actually earned slightly less than subjects in the D-Culture (Experiment 1: C-Culture 206.18 MU vs D-Culture 230.65 MU, Rank-sum  $p=0.001$ ; Experiment 2: C-Culture 160.25 vs D-Culture 211.68, Rank-sum  $p<0.001$ ).

Treatment was randomly assigned at the session level, and subjects who participated in Experiment 1 could not participate in Experiment 2. See Table 1 for descriptive statistics for each experiment.

**Table 1.** Descriptive statistics for each experiment. No subjects dropped out of the experiment or are excluded from analysis.

	# Sessions	# Subjects	Mean Age	Female
Experiment 1	6 (3 C-Culture)	96 (44 C-Culture)	21.7	37%
Experiment 2	10 (6 C-Culture)	122 (66 C-Culture)	21.8	45%

### 3. Experiment 1: Effects of cooperative environment on prosociality

#### 3.1 Experimental design

Experiment 1 tests our prediction that exposure to a laboratory environment where cooperation is a strongly supported equilibrium will lead to more prosociality in subsequent one-shot anonymous games compared to a laboratory environment where cooperation is not an equilibrium. To do so, the RPDs of Stage A described above are followed by a Stage B consisting of four games that are widely used for measuring prosociality (for an overview, see Camerer and Fehr (2002)): a Public Goods Game (PGG), a Trust Game (TG), a Dictator Game (DG), and an Ultimatum Game (UG), played in the listed order. We employ the “strategy method”: subjects enter decisions for each possible player role of each game; then at the end of Stage B, one game is selected at random to actually be played for money, and subjects are randomly assigned to player roles in that game with their action beginning determined by the corresponding choice they indicated earlier. Thus subjects receive no feedback about outcomes between different decisions in Stage B, yet their decisions are still incentivized. Subjects are fully informed about this procedure, as well as the fact that earnings will be translated into cash at the same exchange rate used in Stage A.

In the PGG, players interact in groups of 4. Each player is given an endowment of 100 MU and chooses how much of it to keep for him/herself, and how much to transfer (contribute) to a common project. All MU transferred are multiplied by an efficiency factor of 1.6 and then divided equally among all 4 group members.

In the TG, players interact in pairs: a Trustor and a Trustee, both of whom begin the game with an endowment of 50 MU. The Trustor chooses whether or not to transfer her 50 MU to the Trustee (binary choice). If the transfer is made, the 50 MU are tripled and given to the Trustee. The Trustee chooses how many MU to transfer back to the Trustor (any amount from 0 to 150 MU) should the Trustor make the



transfer. In addition to measuring behavior, we also measure Trustor expectations in this game: players are asked to provide an (unincentivized, as per Peysakhovich and Plagborg-Møller (2012)) guess of the average number of MU transferred back by Trustees.

In the DG, subjects again are assigned to pairs: a Dictator who begins with 100 MU and a Recipient who begins with 0 MU. The Dictator then chooses how much of her endowment to transfer to the Recipient, who is passive and takes no action. Thus there is no role for the Dictator's expectations about the behavior of the Recipient, and the DG offers a pure measure of prosocial preferences.

In the UG, subjects once again interact in pairs: a Proposer and a Responder are given 100 MU to divide between them. The Proposer chooses how many MU to offer to the Responder. The Responder indicates the minimum offer she is willing to accept (minimum acceptable offer, MAO). If the Proposer's offer is greater than or equal to the Responder's MAO, then the transfer occurs and each player is paid accordingly. If, on the other hand, the Proposer's offer is less than the Responder's MAO, the offer is rejected and neither player earns anything.

Thus subjects make six decisions in Stage B. Five of these decisions (all but the UG Responder decision) involve transferring money from oneself to one or more others. These five transfer decisions therefore form our Stage B measures of prosociality in one-shot anonymous settings, which we analyze jointly. To compare transfers across different games and roles, we normalize each decision such that the maximum transfer has a value of 1. Thus a normalized value of 1 is assigned to transferring 100 in the PGG, choosing to transfer as TG Trustor, transferring back 150 as TG Trustee, transferring 100 in the DG, and offering to transfer 100 as UG Proposer. Our choice to take these five decisions as measures of prosociality, and not to include the UG MAO, is motivated by a Principal Component Analysis of Stage B play in Experiment 1, which suggests that the five prosociality decisions track together while the UG MAO does not (see Appendix for details). Further support comes from evidence that these transfer decisions are strongly correlated within an individual, but are not correlated with that individual's UG MAO (Peysakhovich et al., 2014, Yamagishi et al., 2012).

To assess the role heuristics play in any effect Stage A might have on Stage B, we have subjects complete the cognitive reflection test (CRT) after finishing Stage B. The CRT is a set of 3 simple math problems with intuitively compelling but incorrect answers (Frederick, 2005).<sup>11</sup> We use CRT scores as a

---

<sup>11</sup> As the CRT has become a commonly used measure, and prior exposure to the questions may undermine their effectiveness, we use a modified version introduced in (Shenhav et al., 2012), which has the following questions: (Q1) The ages of Mark and Adam add up to 28 years total. Mark is 20 years older than Adam. How many years old is Adam? (Correct Answer: 4, Intuitive Answer: 8); (Q2) If it takes 10 second for 10 printers to print out 10 pages of paper, how many seconds will it take 50 printers to print out 50 pages of paper? (Correct Answer: 10, Intuitive Answer: 50); (Q3) On a loaf of bread, there is a patch of mold. Every day, the patch doubles in size. If it takes 12

proxy for propensity to engage in intuitive thinking (and thus to follow one's heuristic response) vs. stopping to think and thus potentially overriding one's heuristic response.<sup>12</sup> Finally, to assess whether any treatment effect is driven by changing subjects' general affect or mood, subjects in three sessions (N=48, 24 in 2 C-Culture sessions and 24 in 1 D-Culture session) are asked "How would you describe your mood right now?" at the end of the experiment, using a 5-point Likert scale (1- Very bad to 5- Very good).

### 3.2 Results

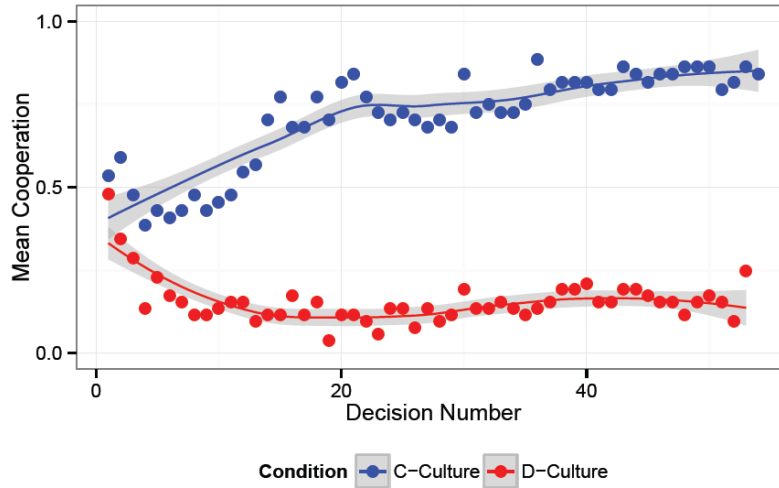
#### 3.2.1 Stage A Cooperation

We first conduct a manipulation check to verify that our Stage A game specifications worked as expected. Indeed, we observe high levels of RPD cooperation in the C-Culture treatment and low levels of RPD cooperation in the D-Culture treatment (Figure 1). Consistent with this observation, a regression (linear probability model) predicting cooperation finds a significant positive effect of C-Culture dummy ( $p < 0.001$ ). A second regression shows that this difference increases as subjects learn over time (significant positive interaction between C-Culture dummy and decision number,  $p < 0.001$ , in a regression including decision number). See Appendix Table A1 for regression table. Thus our Stage A manipulation successfully creates conditions of persistent cooperation (in the C-Culture) or defection (in the D-Culture).

---

days for the patch to cover the entire loaf of bread, how many days would it take for the patch to cover half of the loaf of bread? (Correct Answer: 11, Intuitive Answer: 6)

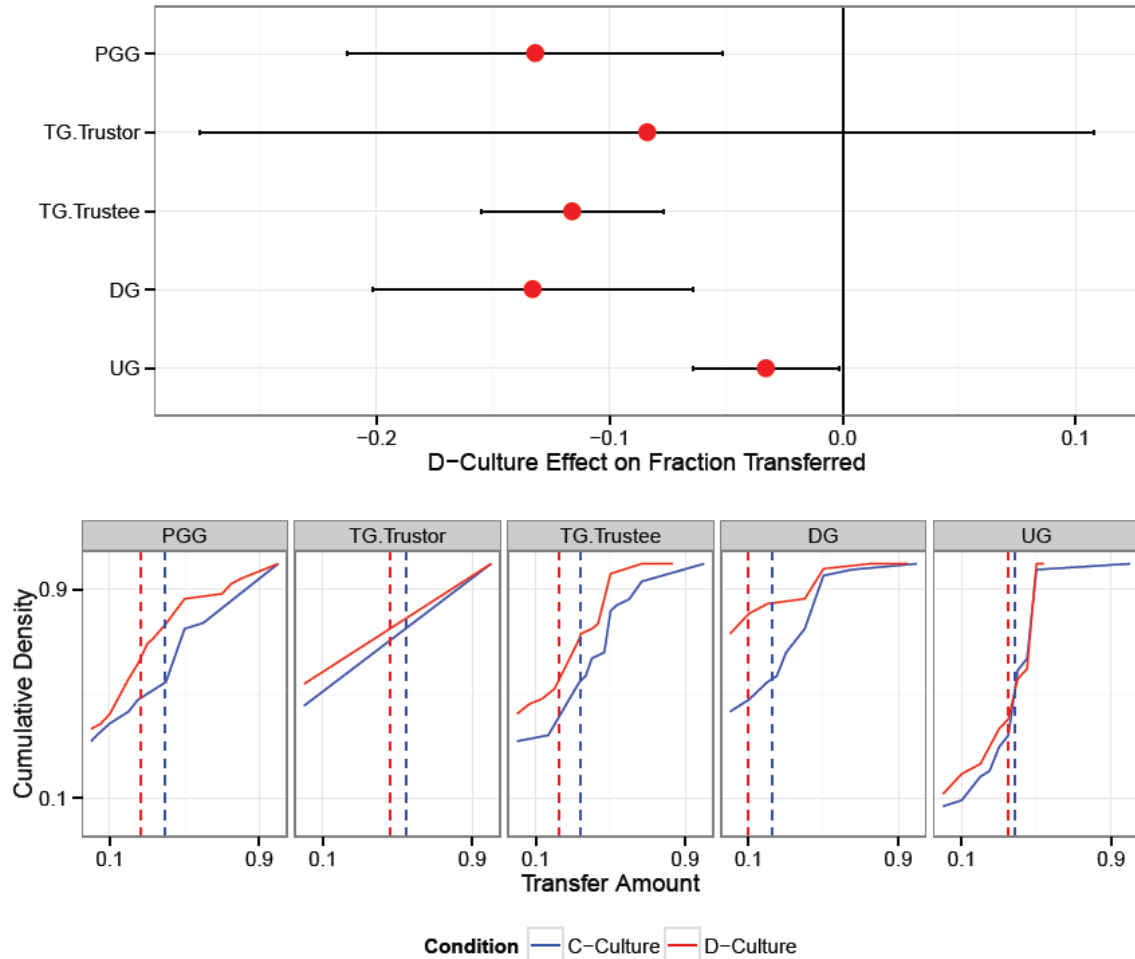
<sup>12</sup> Welsh et al. (2013) present evidence that performance on the CRT correlates with heuristic use to a greater extent in contexts where numerical skill is involved in arriving at the correct answer compared to contexts where numerical skill is not required. Based on these findings, CRT score should successfully tap into heuristic use in our experiment, as the games of our Stage B clearly have a numerical component. To further distinguish between general mathematical ability and reliance on intuition, we include a secondary analysis examining number of intuitive answers provided on the CRT (rather than number of correct answers). We also note that other work has found that CRT correlations (as we use here) and actual experimental manipulations of intuitive processing have similar effects, even in non-numerical domains (e.g. belief in God (Shenhav, et al., 2012)).



**Figure 1:** Fraction of subjects cooperating in each Repeated Prisoner’s Dilemma decision of Stage A in Experiment 1, with locally estimated (LOESS) 95% confidence intervals.

### 3.2.2 Stage B Prosociality

We now turn to our main question of interest: does prosociality in Stage B differ between subjects exposed to the C-Culture and D-Culture treatments in Stage A? As predicted, we see that subjects randomized into the D-Culture treatment transfer substantially less money to others in Stage B than those randomized into the C-Culture (Figure 2;  $p < 0.001$  with or without dummies for decision type, see Appendix Table A3 col 1 and 2). When including interactions between treatment and decision type, we find only a significant interaction with the UG sender dummy ( $p = 0.018$ ; all other interactions  $p > 0.4$ ; see Appendix Table A3 col 3), such that Stage A has less of an effect on UG offers. Nonetheless, Stage A still increases transfers even in the UG: when analyzing each decision separately (Appendix Table A4), we find a significant positive effect in all decisions except TG Trustor. (Although there is a relatively large effect size for TG Trustor, as seen in Figure 2, the standard errors are also larger than for the other decisions because TG Trustors make a binary rather than scalar decision.) Thus we find strong support for our prediction that Stage A repeated game experiences spill over to the strategically different setting of Stage B’s one-shot anonymous games. Exposure to rules that cause subjects (and everyone they interact with) to cooperate or defect almost all of the time alters subsequent prosociality, with strategies that are advantageous over many repeated trials in Stage A being carried over into the novel one-shot decision settings of Stage B.

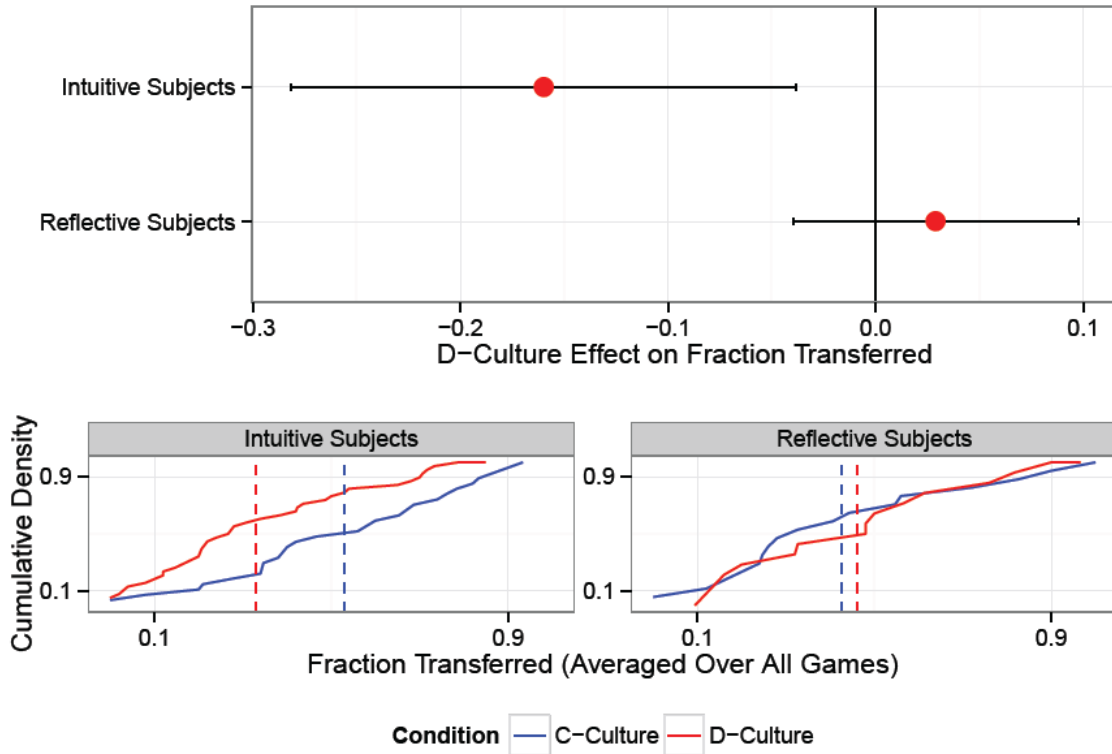


**Figure 2:** Top panel: Mean fraction transferred in D-Culture minus mean fraction transferred in C-Culture for each Stage B decision in Experiment 1, with error bars indicating bootstrapped 95% confidence intervals clustered on session. Bottom panel: Cumulative Distribution Function, CDF, of fraction transferred in each Stage B decision for the C-Culture (blue) and D-Culture (red); mean transfers for each condition are indicated by vertical dashed lines.

### 3.2.3 Role of heuristics in spillovers from Stage A to Stage B

The SHH predicts that the Stage A manipulation will act largely by altering subjects' intuitive responses (rather than their rational, deliberative responses). Thus, Stage A is predicted to have less of an effect on subjects who are better at overruling their heuristic responses when making decisions (and therefore who score better on the CRT). Consistent with this prediction, we find a significant negative interaction between number of correct answers on the CRT and a C-Culture dummy when predicting Stage B

prosociality ( $p=0.028$ ; Appendix Table A5 col 1). This effect is visualized in Figure 3, by comparing the treatment effect among non-deliberative thinkers (those giving 1 or more incorrect CRT answers,  $N=65$ ;  $p=0.010$ ; Appendix Table A5 col 2) and deliberative thinkers (those giving no incorrect CRT answers,  $N=31$ ;  $p=0.413$ ; Appendix Table A5 col 3). Furthermore, although the number of incorrect CRT answers is the standard measure of (non)intuitiveness, subjects might answer incorrectly by giving an answer which is neither the intuitive answer nor the correct answer (perhaps due to limited mathematical ability). Since we are particularly interested in the use of intuition (rather than general inability to do math), we note that we find a similar moderation effect using number of specifically intuitive answers rather than just generally incorrect answers (see Appendix Table A6). Taken together, these results provide evidence that the changes which Stage A causes in Stage B behavior are specifically the result of reshaping intuitive/heuristic responses.



**Figure 3:** Top panel: Mean fraction transferred in D-Culture minus mean fraction transferred in C-Culture averaged on over all five Stage B transfer decisions, for intuitive subjects (those giving 1 or more incorrect CRT answers) and reflective (those giving no incorrect CRT answers) subjects; error bars indicate bootstrapped 95% confidence intervals clustered on session. Bottom panel: CDFs of average fraction transferred across all decisions for the C-Culture (blue) and D-Culture (red); mean transfers for each condition are indicated by vertical dashed lines.

### 3.2.4 *Role of beliefs in spillovers from Stage A to Stage B*

A natural question that arises is whether the effect of Stage A on behavior in Stage B can be explained solely by changes in subjects' beliefs, rather than an actual remodeling of preferences (which are typically assumed to be fixed). The most straightforward way that beliefs could affect behavior is by changing expectations about the decisions of one's Stage B co-player(s). And indeed, we find such an effect: TG Trustors coming from the D-Culture expect a significantly smaller back-transfer from Trustees than those from the C-Culture; 33 MU vs 48 MU,  $p < .001$ ; Appendix Table A4 col 6). Changing expectations about the behavior of one's partner cannot, however, explain all of our results: we find treatment effects in decisions that are not influenced by such expectations. Trustee decisions in the TG are explicitly conditioned on their Trustor's behavior (the Trustee's decision is only implemented if the Trustor chooses to transfer); therefore, expectations of Trustor behavior have little room to influence Trustee behavior. Even clearer is the DG, where the recipient makes no decision, and thus there is no room whatsoever for expectations about the partner's decision. Yet we find large and statistically significant treatment effects on both of these decisions (TG Trustee: 33.7% vs 22.1% transferred,  $p < 0.001$ , Appendix Table A4 col 4; DG: 22.7% vs 9.5% transferred,  $p < 0.01$ , Appendix Table A4 col 1). These effects cannot be explained by beliefs about the partner's choices.

There is, however, a somewhat more subtle way that beliefs might influence behavior, as captured by social preference models of type-based reciprocity (e.g. Levine (1998)). These models suggest that one's desire to help another person depends on one's beliefs about how helpful that other person is. Thus, in the DG, the dictator's beliefs about the altruism of their recipient (e.g. what the recipient would have done had he or she been the dictator) might influence the dictator's behavior, even though the recipient makes no actual decision in the game. When making their DG decisions in Stage B, subjects do not know anything about their recipient's history of play in Stage A. They do, however, know that their recipient also participated in Stage A. Thus dictators might draw inferences about their recipient's type based on prior experience with their various Stage A partners, allowing Stage A to influence DG giving in Stage B via beliefs about type (although such inferences would likely be incorrect, since random assignment ensures that differences in Stage A play across conditions are purely a result of the rules of the game, rather than differences in the distribution of types across conditions).

Here we present two pieces of evidence that speak against this possibility. The first comes from comparing the treatment effect on the DG versus the TG Trustee. In the DG, there is ambiguity about how one's partner would have acted (leaving room for type-based beliefs to affect behavior). This is substantially less true in the TG, however, because the Trustee's decision is only implemented if the

Trustor transfers money. Therefore, there is less ambiguity about the partner's type when choosing how much to return: Trustees know that their decision will only be implemented if their partner made the "nice" choice of transferring. This implies that if the treatment effect is driven solely by changing expectations about the partner's type, we should see a substantially larger treatment effect for the DG (where there is more type ambiguity) than the Trustee. Contrary to this prediction, however, we observe very similar effect sizes in the two games: when comparing the D-Culture and the C-Culture, 11.6% less of the stake is transferred by Trustees and 13.2% less is transferred by Dictators, a non-significant difference (Interaction between treatment dummy and Trustee dummy,  $p=0.68$ ; Appendix Table A3 col 3). Thus, reducing the amount of ambiguity regarding the partner's type does not reduce the treatment effect.

Our second piece of evidence comes from an additional Supplemental Experiment designed to minimize the inferences subjects draw regarding the type of their Stage B partner based on their Stage A interactions. Specifically, our Supplemental Experiment's Stage B consists only of a single DG where the recipient does not participate in Stage A; the DG recipients are completely passive, and the dictators are informed of this fact very prominently. Here, the play of others in Stage A should have no bearing on one's expectations regarding the Stage B recipient. We recruit 237 American subjects using the (non-student based) online labor market Amazon Mechanical Turk (Amir et al., 2012, Horton et al., 2011). We then have them play an adapted version of our Stage A RPD followed by a dictator game in which the recipient has not played the RPD, and receives no payment other than what is given to them in the DG. The results of this experiment closely replicate the DG results observed in Experiment 1, with nearly twice as much DG giving after the C-Culture compared to the D-Culture (27.4% vs 16.5% transferred, Rank-sum  $p<0.001$ ). See Appendix Section 4 for details. Thus reducing the grounds for using Stage A to make inferences about the type of the DG recipient in Stage B does not appreciably undermine the treatment effect.

In sum, while we cannot completely rule out the possibility that our treatment effect is driven by beliefs about the distribution of types, we provide evidence that is inconsistent with this possibility, and instead points to actual remodeling of preferences (e.g. internalized notions of appropriate behavior).

### *3.2.5 Role of mood in spillovers from Stage A to Stage B*

Finally, we note that our treatment effect does not appear to be driven by mood. Most simply, we find little difference in our 1-to-5 mood measure between the C-Culture (mean=3.45, std=0.588) and the D-Culture (mean=3.29, std=0.690; Rank-sum,  $p=0.38$ ) in the 3 sessions of Experiment 1 in which mood was measured after Stage B. (We also find no significant effect of condition on mood in the 9 sessions of

Experiment 2 where mood was measured: C-Culture, mean=3.43, std=0.095; D-Culture, mean=3.36, std=0.106; Rank-sum, p=0.46). Further evidence comes from replicating our main analysis using only these 3 sessions (Appendix Table A7). If the effect of treatment is driven by Stage A altering subjects' mood, then the treatment coefficient will be smaller when controlling for mood. On the contrary, however, we find that the coefficient on treatment *increases* in magnitude when including mood in the regressions (C-Culture dummy: coeff=0.128 without mood, coeff=0.139 with mood). Thus we do not find evidence that our treatment effect is driven by mood. We do note, however, that mood was measured after Stage B decisions were made, and thus is not a pure measure of Stage A's effect on mood. Thus further investigation of the role of mood in our treatment effect is a worthwhile direction for future work.

#### 4. Experiment 2: Effects of cooperative environment on norm enforcement

##### 4.1 Experimental design

A key element of norms is not just the desire to act in a particular way, but also the willingness to sanction those who do not act accordingly. In Experiment 2, we therefore ask whether Stage A also affects subjects' punishment of selfishness (i.e. willingness to pay to reduce the payoffs of non-cooperators). To do so, we follow Stage A with a battery of one-shot punishment games, played using the strategy method with no feedback: a DG with third party punishment (3PDG), a Prisoner's Dilemma with third party punishment and reward (3PPD), and a PGG with punishment (PGP). In each of these punishment decisions, players face a different type of dilemma from Experiment 1: they have the chance to decrease the payoff of one or more others at a cost to themselves, based on the others' behavior.

In the 3PDG game, subjects are matched in groups of three. One subject is assigned to be the dictator and unilaterally decides on a split of 100 MU between herself and a second subject assigned to be the recipient. A third subject, the sanctioner, is given an endowment of 100 MU and indicates how many MU (up to 20) she would spend to reduce the dictator's payoff, depending on the dictator's chosen division.<sup>13</sup> Each MU spent by the third party reduces the dictator's payoff by 5 MU.

In the 3PPD game, subjects are matched into groups of four. Two subjects are selected to be PD players and play a one-shot binary-choice PD with each other using the following payoff matrix (in MU):

	C	D
C	80,80	0,120
D	120,0	20,20

<sup>13</sup> The dictator had the option of transferring 0, 10, 20, 30, 40, or 50 MU. The sanctioner indicated the number of MU to spend on reducing the dictator's payoff for each of these possible transfer options.



The other two subjects are sanctioners and are each given an endowment of 100 MU. The sanctioners choose how much (up to a maximum of 20 MU) to spend on reducing one of the two PD players' payoff, based on that player's decision in the PD. We also use this game to address concerns that have been raised about third-party punishment experiments where punishment was the only option (Pedersen et al., 2013): we allow third parties to reward as well as punish the PD players. Each MU spent by the third party on punishing reduces the PD player's payoff by 5 MU, and each MU spent on rewarding increases the PD player's payoff by 5 MU.

In the PGP game, subjects are matched into groups of four and play the same PGG as in Stage B of Experiment 1 (100 MU endowment, efficiency factor of 1.6). Unlike Experiment 1, however, after making a PGG decision each subject has the opportunity to pay up to 20 MU to sanction other group members based on contribution amount.<sup>14</sup> Each MU spent reduces the sanctioned player's payoff by 5 MU.

In all three games, in addition to punishment decisions, subjects are asked to rate how "socially inappropriate" each possible punishee behavior was, using a 7-point Likert scale.

## *4.2 Results*

### *4.2.1 Stage A Cooperation*

As in Experiment 1, we effectively induce substantially more Stage A cooperation in the C-Culture treatment compared to the D-Culture treatment (58% C vs 22% C,  $p < 0.001$ ; See Appendix Table A11).

### *4.2.2 Stage B Third Party Punishment*

We begin by considering the effect of Stage A on third party punishment in Stage B (Figure 4). In the third party punishment games (3PDG & 3PPD), punishment is impartial (as the punisher is not affected by the behavior of the punishee) and thus reflects pure norm enforcement (Fehr and Fischbacher, 2004b).

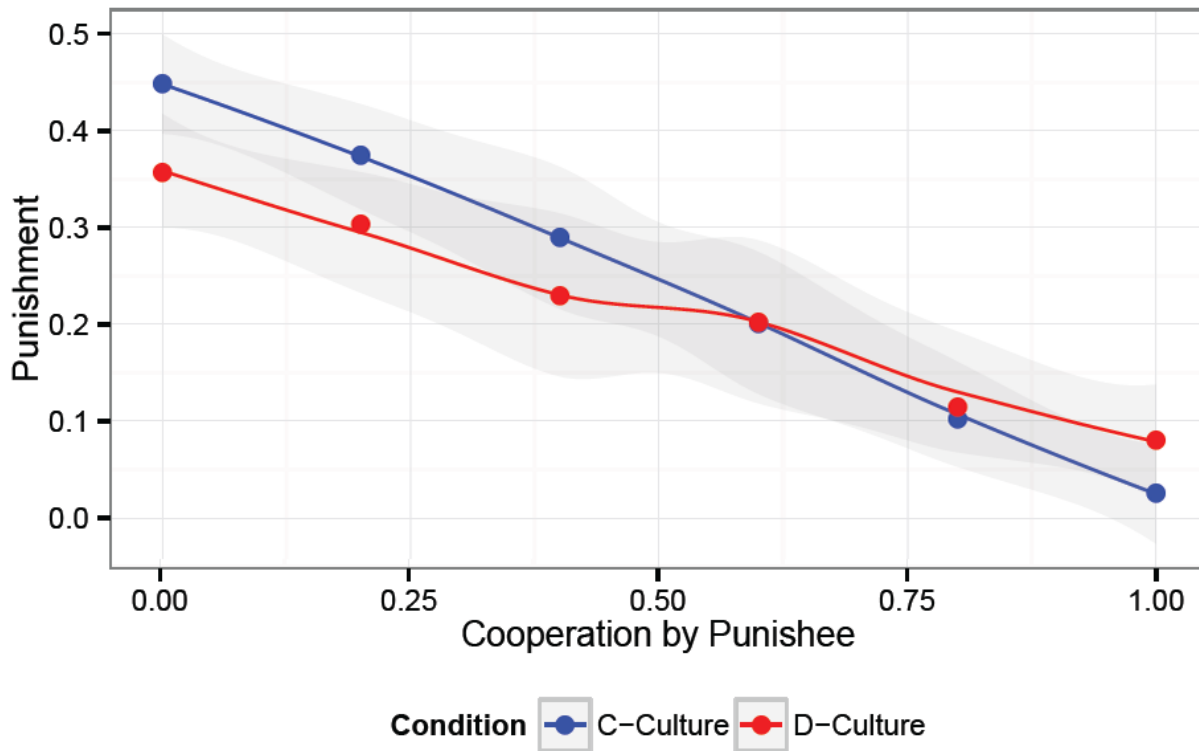
---

<sup>14</sup> In the PGP contribution phase, subjects can choose one of the following contribution amounts: 0, 25, 50, 75, 100. In the PGP punishment phase, subjects are paired up and indicate how many MU to spend punishing the person with whom they are paired, for each possible partner contribution amount. We do not offer a reward option in the PGP, but based on the similarity of results between of our third-party punishment games with and without reward, we think it is unlikely that adding a reward option to the PGP would change our results substantially.

To test whether Stage A affected norm enforcement in Stage B, we follow the approach of Experiment 1 and analyze the data from the 3PPD and the 3PDG together in a single regression. To make cooperation in the DG and PD comparable, we normalize such that maximum selfishness (transferring nothing in DG, playing D in PD) has a value of 0 and maximum prosociality (transferring 50 in the DG, playing C in PD) has a value of 1. Regressing amount of punishment on the punishee's level of cooperation, a dummy for C-Culture, and an interaction between punishee's level of cooperation and the C-Culture dummy (Appendix Table A12 col 1) shows a significant negative main effect of punishee's cooperation (more cooperative actions were punished less,  $p < 0.001$ ), a significant positive main effect of C-Culture (a maximally selfish action was punished more in the C-Culture than in the D-Culture,  $p = 0.021$ ) and a significant negative interaction (for a given increase in punishee's cooperation, punishment decreased more in the C-Culture than in the D-Culture; that is, punishment was more selectively targeted at selfish behavior in the C-Culture,  $p = 0.006$ ).

These results are robust to including fixed effects for game type (Appendix Table A12 col 2) as well as interactions between game type and treatment (Appendix Table A12 col 3), or to analyzing the two games separately (Appendix Figure A3 and Table A12 col 4 and 5). We also find that subjects' ratings of the inappropriateness of DG sending and PD cooperation decisions are affected by the treatment in a qualitatively similar way (i.e. positive C-Culture dummy coefficient, negative interaction between C-Culture and Punishee Prosociality), but that these self-report measures are less sensitive than actual punishment and statistical significance is not achieved (Appendix Table A13); and that the treatment does not have a significant effect on rewarding in the 3PPD (Appendix Table A12 col 6).

Importantly, we find significant treatment effects on third party punishment among intuitive subjects but not among reflective subjects (Appendix Table A14), replicating the pattern seen for prosociality in Experiment 1.



**Figure 4:** Third-party punishment in the 3PDG and 3PPD games of Stage B in Experiment 2, as a function of the punishee’s cooperativeness (fraction of punishment endowment spent, averaged across the 3PDG and 3PPD games), with LOESS 95% confidence intervals. Punishee’s cooperativeness for the 3PDG is normalized such that giving nothing corresponds to cooperativeness 0 and giving 50% corresponds to cooperativeness of 1; and for the 3PPD such that defecting corresponds to cooperativeness 0 and cooperating to cooperativeness 1.

These results demonstrate that Stage A altered subjects’ conception of cooperation as a norm to be enforced as well as followed, rather than just priming subjects to themselves give or not: subjects in the C-Culture condition not only behaved more cooperatively in Experiment 1 but were also more willing to pay to punish non-cooperators in Experiment 2.

#### 4.2.3 Stage B Public Goods Game with Punishment

We now turn to the effect of Stage A on punishment in the Public Goods Game. Unlike the 3PDG and 3PPD games, punishment in the public goods game is not wholly impartial: non-contributors have a negative impact on the earnings of the potential punisher as well as on other group members. Thus motives other than norm enforcement, such as retaliation or desire to out-earn other group members, may

also drive punishment in this setting (Bolton and Ockenfels, 2000, Ellingsen, et al., 2012, Espín et al., 2012, Herrmann, et al., 2008, Rand and Nowak, 2011).

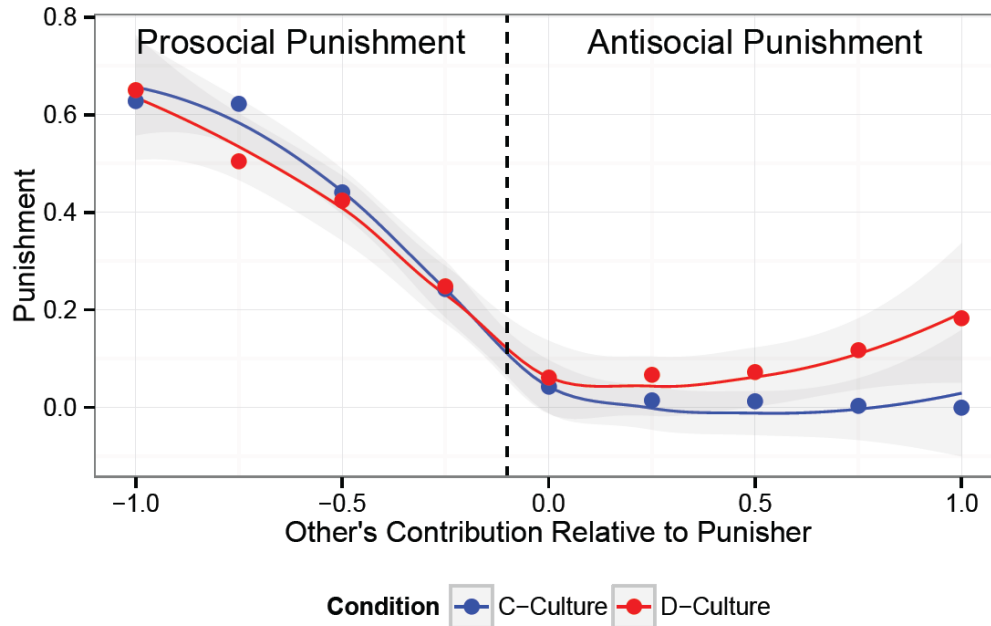
Perhaps as a result of these mixed motives, cross-cultural experiments using the public goods game (e.g. Herrmann, et al. (2008)) have found relatively little variation in prosocial punishment (punishment of those contributing less than the punisher, PSP). The level of “antisocial punishment” (punishment of those contributing as much or more than the punisher, ASP), conversely, has been shown to vary dramatically across cultures. Thus cross-cultural effects on PGG punishment vary based on the difference between the punisher’s contribution and the punishee’s contribution.

We therefore consider Stage A effects on PGG punishment based on this difference in punisher and punishee contribution levels (Figure 5).

First we consider PSP of those who contribute less than the punisher (left half of Figure 5). Consistent with previous cross-cultural results (Herrmann, et al., 2008), we find no significant difference in overall PSP between our C-Culture and D-Culture treatments (Appendix Table A15 col 1,  $p=0.72$ ), and no interaction between treatment and the difference in punisher’s and punishee’s contribution amount (Appendix Table A15 col 2,  $p=0.65$ ).

Next we turn to ASP targeted at those who contributed as much or more than the punisher (right half of Figure 5). Our results are again consistent with (Herrmann, et al., 2008)’s cross-cultural work, where such punishment is rare among Western college students but occurs frequently in countries with weaker institutions: We find virtually no antisocial punishment in the C-Culture treatment, but some begins to appear in the D-Culture treatment. There is a trend in the direction of less overall ASP in the C-Culture (Appendix Table A15 col 3,  $p=0.16$ ), and there is a marginally significant interaction between the C-Culture dummy and contribution difference (Appendix Table A15 col 4,  $p=0.098$ ) such that in the D-Culture, subjects are somewhat more likely to punish as the punishee’s contribution *increases*.

These results provide further evidence that our Stage A manipulation is altering the norms applied in Stage B, and hint that cross-cultural differences in ASP (Ellingsen, et al., 2012, Herrmann, et al., 2008) may in part be explained by norms developed in an environment where cooperative behavior is not an advantageous strategy (a possibility which deserves further exploration in future research).



**Figure 5:** Punishment in the PGP of Stage B in Experiment 2, as a function of the difference between the punishee’s normalized PGG contribution and the punisher’s normalized PGG contribution (with LOESS 95% confidence intervals). Negative x-axis values correspond to prosocial punishment (PSP, where the punisher contributes more than the punishee) whereas positive x-axis values correspond to antisocial punishment (ASP, where the punishee contributes more than the punisher).

## 5. Beyond play in economic games

Finally, we ask whether the effects of Stage A extend beyond play in economic games. To do so, we employ the World Values Survey (WVS), one of the most widely used instruments for the measure of differences in cultural norms (Gächter et al., 2004, Putnam, 2000). A version of the World Values Survey question measuring generalized trust was included in the post-experimental survey in each of our experiments<sup>15</sup>: subjects were asked “How much do you agree with the statement: ‘most people can be trusted.’?” and indicated their response using a Likert scale ranging from 1 [Strongly disagree] to 5 [Strongly agree]. This specific continuous implementation is used by Peysakhovich, et al. (2014), who

<sup>15</sup> Certain subjects’ questionnaires did not include this question, but not in a systemic way: three sessions of Experiment 1 contained the trust question (48 observations), all but 1 session from the Experiment 2 contained the trust question (109 observations), and all 237 observations from the Supplemental Experiment described in the Appendix included the trust question.

find that agreement with the statement correlates strongly with prosociality in the one-shot anonymous PGG and DG. The correlation with DG giving suggests that this question measures more than simple beliefs about the behavior of others that is implied by the word “trust”, but instead is a more general measure of the responder’s prosociality. Additional evidence for this broader interpretation of the WVS trust question comes from Glaeser et al. (2000), who find that it predicts trustworthiness much better than trust. Thus, it seems likely that the Stage A effects on prosocial norms observed in the game behavior of Experiments 1 and 2 may also translate into effects on this generalized trust measure.

Consistent with our results regarding play in economic games, the Stage A treatment does indeed have a substantial effect on our WVS trust measure. Subjects randomized into the D-Culture report significantly lower levels of generalized trust than those randomized into the C-Culture, both overall (Appendix Table A16 col 1,  $p < 0.001$ ) and among the lab (Appendix Table A16 col 3,  $p = 0.040$ ) and online subjects (Appendix Table A16 col 4,  $p = 0.001$ ) separately. Furthermore, we find evidence of a similar CRT moderation effect to that observed with the games: Stage A significantly changes WVS trust among intuitive subjects, but not among reflective subjects (see Appendix Tale A17).

Thus, the effect of our Stage A manipulation is not restricted to economic games; we also recreate variation in a standard survey instrument used to measure culture through random assignment to an interaction environment in the lab. With this measure, the heuristics developed in Stage A are misapplied not to behavior in one-shot anonymous games, but rather to a survey question that combines assessments about how trustworthy others are with one’s own level of prosociality.

## **6. Concluding discussion**

Here we have shown that externally imposed interaction rules can dramatically alter subjects’ internalized norms: immersing subjects in an environment that incentivizes cooperation or defection for less than 20 minutes leads to large differences in subsequent one-shot anonymous prosociality and sanctioning, as well as generalized trust. Furthermore, these effects are driven by subjects that rely on heuristics, and virtually disappear among highly deliberative subjects. Our findings demonstrate the important part that spillover effects can play in shaping behavior in one-shot anonymous settings. Furthermore, we shed light on the key role played by intuitive and heuristics in this internalization of cooperative norms (as proposed by the Social Heuristics Hypothesis (Rand, et al., 2014c)). In doing so, we help to explain why people often cooperate in one-shot anonymous contexts, and why such behavior might vary across organizations and cultures. More broadly, our results demonstrate that laboratory experiments can be a powerful tool for studying culture and the interplay between social environments and internalized norms.

Our results cannot be explained by changes in beliefs regarding the decisions of co-players, as we see treatment effects in the Dictator Game and the Trust Game Trustee. We also present evidence that our treatment effects are not easily explained by changes in beliefs regarding the co-player's type (see Section 3.2.4), although we cannot definitively rule out this possibility. We argue that our results therefore suggest that social *preferences* may be remodeled on relatively short time-scales (20 minutes or less of RPD play).<sup>16</sup> This suggestion runs counter to the standard assumption in economic models that preferences are fixed, and points to the importance of considering both where social preferences come from and how they change. These questions are not addressed by most social preferences models, which take preferences as given (Bolton and Ockenfels, 2000, Charness and Rabin, 2002, Fehr and Schmidt, 1999, Levine, 1998, Rabin, 1993). In particular, the evidence we present for the role of heuristics in this remodeling of preferences suggests that adding a dual process perspective to social preference models (e.g. Dreber et al. (2014), Loewenstein and O'Donoghue (2004)) is important when exploring preference change. Evolutionary game theoretic models *do* often try to shed light on the issues of origin and change of prosocial behavior (e.g. Alger and Weibull (2010), Bowles (2001), Dekel et al. (2007), Rand and Nowak (2012), Rand et al. (2013), Samuelson (2001)), but often without explicitly connecting to specific models of social preferences. Extending social preference models to incorporate heuristic processing and preference change, and more directly linking evolutionary models to models of social preferences are important directions for future theoretical work.

It is also interesting to consider how our results relate to the economic theory literature on habit formation, growing out of the seminal work of Becker and Murphy (1988). The Becker-Murphy model defines habit formation as adaptation to levels of consumption of a particular good: the amount of the good you desire to consume in the current period is increasing in the amount you consumed in the previous period. To apply this model to our experiments, think of PD cooperation as the good being consumed in Stage A. In order to explain the spillover over effects we observe in Experiment 1, we would have to assume that subjects think of RPD cooperation in Stage A and money transfers in the one-shot games of Stage B as instances of the same 'prosociality' good: subjects in the C-Culture consume more prosociality in the Stage A RPD than those in the D-Culture, and therefore are more inclined to consume prosociality in the anonymous one-shot games of Stage B. This assumption may be reasonable, given evidence that an individual's play correlates strongly across different prosociality games, suggesting that a common preference drives behavior across these games (Peysakhovich, et al., 2014). To explain Experiment 2, conversely, we would have to assume that subjects experience Stage A RPD cooperation

---

<sup>16</sup> Convergent empirical evidence for such remodeling comes from studies showing that short-run factors such as recent experiences or organizational memberships can affect internalized behaviors related to trust and cooperation (e.g. Alesina and La Ferrara (2002), Bellows and Miguel (2009), Fisman and Khanna (1999)).

and Stage B third party punishment of selfishness as being the same good. This assumption is contradicted, however, by evidence that an individual's play in cooperation games is *not* predictive of their third party punishment behavior (Peysakhovich, et al., 2014). This suggests that cooperation and punishment are psychologically distinct goods, and that our results are driven by norm internalization rather than by pure habituation to cooperation. Building models that connect the economic literature on habit formation to the concept of norms and norm internalization is an important avenue for future work.

While we demonstrate change in behavior following relatively brief exposure to cooperative or non-cooperative environments, we do not systematically vary treatment durations. Thus we cannot estimate how different lengths of exposure to the treatment translate into different levels of change in subsequent behavior. Creating such a 'dose-response' curve is an important direction for future experimental work. So too is exploring the extent to which the Stage A treatment generalizes across contexts. We provide some evidence of generalization, in that RPD behavior generalizes to other prosociality games, to punishment games, and to the non-game survey measure of trust from the World Values Survey. We do not contend, however, that our Stage A treatment erases a lifetime of previous experience and cultural context, or irrevocably change our subjects' preferences. Rather, our experiments should be seen as a proof-of-concept that the behaviors cultivated by particular environments travel beyond these situations, and that cooperation includes relatively malleable components. Future work should evaluate how far beyond the experimental context our effect extends, as well as how our artificial Stage A can be translated into more contextualized settings (which may lead to even broader generalization).

A related question concerns the persistence of the effect we observe. Some insight comes from the findings of Duffy and Ochs (2009) and Fréchette and Yuksel (2013), who examine play in a series of RPD games: in a first phase subjects played RPDs where cooperation is an equilibrium, then in a second stage the game specification is modified such that cooperation is not an equilibrium.<sup>17</sup> Both studies find a large reduction in cooperation when switching from the first stage to the second. Consistent with our findings, however, there is some indication of spillover: cooperation in the first decision of the second stage is higher than in later periods, and a process of decay occurs before aggregate behavior stabilizes at an extremely low level. This rapid decay indicates that subjects quickly adapt to their new environment. However, this does not mean that these spillover effects are trivial, for two reasons. First, in real-world applications, environmental conditions are themselves typically highly persistent: people interact every day under an environment that favors either cooperation or defection. Thus they "receive treatment"

---

<sup>17</sup> This switch is accomplished by switching from fixed matching to random matching in the case Duffy and Ochs (2009) and by changing the stage game payoffs in the case of Fréchette and Yuksel (2013).



constantly, supporting a continual spillover of norms into the subset of interaction that are one-shot and anonymous. Second, many situations of economic interest have multiple equilibria (unlike our Stage B games or those of Duffy and Ochs (2009) and Fréchet and Yuksel (2013)). Thus even short-lived spillover effects can have long lasting impacts by changing the initial conditions and shifting between basins of attraction of different equilibria, as demonstrated in the context of weak-link coordination games by Brandts and Cooper (2006). Quantifying the persistence of our Stage A effect, and its consequences in settings with multiple equilibria, is an important direction for future experiments.

In our experiments, we modeled environments that make cooperation advantageous or disadvantageous using the framework of repeated games. We used repetition because it is a well-established paradigm that creates future consequences for today's actions. Critically, we are not claiming that repetition *itself* varies across cultures. Instead, we are arguing that variation in whether institutional incentives support cooperative equilibria (*modeled* in our experiments by repetition) leads to variation in internalized norms. Other work on path-dependent preferences suggests that the specific form taken by these institutional incentives also matters. For example, Bohnet and Baytelman (2007) found that adding communication, punishment, or fixed repetition to the TG increases trust and trustworthiness almost exclusively due to changes in expectations (and if anything, these cooperation-inducing institutions 'crowd out' intrinsic preference-based trust, rather than positively influencing social preferences as in our experiment); Bohnet and Huck (2004) found no differences in 1-shot TG play based on whether subjects previously played a series of stranger-matched (i.e. 1-shot) TGs or repeated play, and find some evidence that stranger-matched play with a reputation system actually decreases subsequent trust; in a sequential step-level PGG, Cooper and Stockman (2011) found relatively little effect of prior experience under rules that emphasized different kinds of equity concerns, and whatever effects they did find were very short lived; and Herz and Taubinsky (2013) find that prior experience in markets with substantial competitive pressure results in a large and persistent decrease in fairness concerns in the UG. Given this heterogeneity, future work directly comparing the effect of different environmental structures that incentivize cooperative behavior (e.g. institutions such as markets, democratic governance or centralized reward and punishment) is needed.

We also do not claim that the arrow from institutionally created environments to preferences goes only in one direction – the co-evolution of norms and interaction environments (e.g. institutions) involves a feedback loop between the two. In situations where institutions (e.g. Stage A rules) work by creating cooperative equilibria, non-cooperative equilibria often also exist.<sup>18</sup> Thus the effectiveness of an

---

<sup>18</sup> Note that here we discuss multiple equilibria in the context of Stage A, rather than above where we discuss possible effects of Stage A on Stage B interactions that involve multiple equilibria.

institution is in part determined by the norms of those whom the institution governs. Multiple equilibria can cause norm persistence (rather than the malleability seen in our studies), as demonstrated for example by Nunn and Wantchekon (2011), who show that individuals whose ancestors come from areas in Africa that had been more affected by the slave trade continue to have weaker levels of interpersonal trust today.<sup>19</sup> Future laboratory studies should examine how institutional rules and baseline norms interact. An important part of such work will involve cross-cultural studies, where our Stage A manipulation is applied to groups with differing baseline levels of cooperation. Such studies will shed light on whether the treatment effect we observe is driven by the C-Culture increasing cooperation, the D-Culture decreasing cooperation, or both. Based on the SHH, we predict that baseline behavior in one-shot anonymous games will resemble Stage B of the C-Culture treatment in places with strong institutions, such as the United States; and will resemble Stage B of the D-Culture treatment in places with weak institutions (e.g. Gächter, et al. (2010)).

Research in the social and behavioral sciences is increasingly focusing not just on generating insights into the basic science underlying human behavior, but also on applying these insights to institutions, markets, incentives and organizations outside the laboratory (Fudenberg and Peysakhovich, 2014, Gerber and Rogers, 2009, Gneezy and List, 2006, Rand et al., 2014d, Roth, 2002, Thaler and Benartzi, 2004, Thaler and Sunstein, 2008, Yoeli et al., 2013). The experiments presented here have clear practical implications for building cooperative cultures in organizations. They suggest that cross-organizational differences in culture are in large part determined by cross-organizational differences in which behaviors (cooperative or non-cooperative) are rewarded. These differences in optimal behavior are strongly influenced by the organization's institutions (i.e. how incentives are structured, the focus of the current paper). However, the way in which people are selected to join the organization also has an important effect on organizational culture: as discussed above, the right incentives can create cooperative equilibria, but groups can still coordinate on non-cooperative equilibria if many of the people entering the group are initially non-cooperative. Thus it is ideal for organizations to both design effective institutions, and (at least to some extent) to avoid hiring non-cooperative individuals.

Taken together, the experiments we present here demonstrate the power of previous experience for shaping our behavior in one-shot anonymous settings. They also open the door for a wide array of possible applications for organizations interested in increasing the incidence of cooperative behavior. Prosociality induced by environmental constraints can have a dramatic influence on behavior even in situations where these rules do not apply.

---

<sup>19</sup> Existing theories for such persistence point to channels such as parents 'instilling' values in their children (e.g. Bisin and Verdier (2000), Tabellini (2007)).

## References

- A. Alesina, La Ferrara, E. 2002. Who trusts others? *Journal of public economics*. **85**(2) 207-234.
- I. Alger, Weibull, J.W. 2010. Kinship, Incentives, and Evolution. *The American Economic Review*. **100**(4) 1725-1758.
- O. Amir, Rand, D.G., Gal, Y.K. 2012. Economic Games on the Internet: The Effect of \$1 Stakes. *PLoS ONE*. **7**(2) e31461.
- G.S. Becker, Murphy, K.M. 1988. A theory of rational addiction. *The Journal of Political Economy* 675-700.
- J. Bellows, Miguel, E. 2009. War and local collective action in Sierra Leone. *Journal of Public Economics*. **93**(11) 1144-1157.
- A. Bisin, Verdier, T. 2000. " Beyond The Melting Pot": Cultural Transmission, Marriage, And The Evolution Of Ethnic And Religious Traits. *Quarterly Journal of Economics* 955-988.
- M. Blonski, Ockenfels, P., Spagnolo, G. 2011. Equilibrium Selection in the Repeated Prisoner's Dilemma: Axiomatic Approach and Experimental Evidence. *American Economic Journal: Microeconomics*. **3**(3) 164-192.
- I. Bohnet, Baytelman, Y. 2007. Institutions and Trust Implications for Preferences, Beliefs and Behavior. *Rationality and Society*. **19**(1) 99-135.
- I. Bohnet, Huck, S. 2004. Repetition and reputation: Implications for trust and trustworthiness when institutions change. *American Economic Review* 362-366.
- G.E. Bolton, Ockenfels, A. 2000. ERC: A Theory of Equity, Reciprocity, and Competition. *The American Economic Review*. **90**(1) 166-193.
- S. Bowles. 1998. Endogenous preferences: The cultural consequences of markets and other economic institutions. *Journal of Economic Literature* 75-111.
- S. Bowles. 2001. Individual interactions, group conflicts, and the evolution of preferences. *Social dynamics*. **155** 190.
- S. Bowles, Gintis, H. 2002. *Prosocial emotions*.
- S. Bowles, Gintis, H. 2003. Origins of human cooperation. *Genetic and cultural evolution of cooperation* 429-443.
- R. Boyd, Richerson, P.J. 2009. Culture and the evolution of human cooperation. *Philosophical Transactions of the Royal Society B: Biological Sciences*. **364**(1533) 3281-3288.
- J. Brandts, Cooper, D.J. 2006. A change would do you good.... An experimental study on how to overcome coordination failure in organizations. *The American Economic Review* 669-693.
- C.F. Camerer, Fehr, E. 2002. Measuring social norms and preferences using experimental games: A guide for social scientists.
- A.W. Cappelen, Moene, K.O., Sørensen, E.Ø., Tungodden, B. 2013. Needs Versus Entitlements - An International Fairness Experiment. *Journal of the European Economic Association*. **11**(3) 574-598.
- G. Charness, Rabin, M. 2002. Understanding Social Preferences with Simple Tests. *Quarterly Journal of Economics*. **117**(3) 817-869.
- M. Chudek, Henrich, J. 2011. Culture gene coevolution, norm-psychology and the emergence of human prosociality. *Trends in cognitive sciences*. **15**(5) 218-226.
- J. Cone, Rand, D.G. 2014. Time Pressure Increases Cooperation in Competitively Framed Social Dilemmas. *PLoS ONE*. **9**(12) e115756.
- D.J. Cooper, Stockman, C.K. 2011. History dependence and the formation of social preferences: an experimental study. *Economic Inquiry*. **49**(2) 540-563.
- G. Cornelissen, Dewitte, S., Warlop, L. 2011. Are Social Value Orientations Expressed Automatically? Decision Making in the Dictator Game. *Personality and Social Psychology Bulletin*. **37**(8) 1080-1090.
- P. Dal Bó. 2005. Cooperation under the shadow of the future: experimental evidence from infinitely repeated games. *American Economic Review*. **95**(5) 1591-1604.
- P. Dal Bó, Fréchette, G.R. 2011. The Evolution of Cooperation in Infinitely Repeated Games: Experimental Evidence. *American Economic Review*. **101**(1) 411-429.
- E. Dekel, Ely, J.C., Yilankaya, O. 2007. Evolution of Preferences. *The Review of Economic Studies*. **74**(3) 685-704.
- A. Dreber, Fudenberg, D., Levine, D.K., Rand, D.G. 2014. Altruism and Self-Control. Available at SSRN.
- A. Dreber, Rand, D.G., Fudenberg, D., Nowak, M.A. 2008. Winners don't punish. *Nature*. **452**(7185) 348-351.

- J. Duffy, Ochs, J. 2009. Cooperative behavior and the frequency of social interaction. *Games and Economic Behavior*. **66**(2) 785-812.
- T. Ellingsen, Herrmann, B., Nowak, M.A., Rand, D.G., Tarnita, C.E. 2012. Civic Capital in Two Cultures: The Nature of Cooperation in Romania and USA. *Availabe at SSRN*: <http://ssrn.com/abstract=2179575>.
- S. Epstein, Pacini, R., Denes-Raj, V., Heier, H. 1996. Individual differences in intuitive-experiential and analytical-rational thinking styles. *Journal of Personality and Social Psychology*. **71**(2) 390-405.
- A.M. Espín, Brañas-Garza, P., Herrmann, B., Gamella, J.F. 2012. Patient and impatient punishers of free-riders. *Proceedings of the Royal Society B: Biological Sciences*. **279**(1749) 4923-4928.
- E. Fehr, Fischbacher, U. 2004a. Social norms and human cooperation. *Trends in Cognitive Sciences*. **8**(4) 185-190.
- E. Fehr, Fischbacher, U. 2004b. Third-party punishment and social norms. *Evolution and Human Behavior*. **25**(2) 63-87.
- E. Fehr, Gächter, S. 2000. Cooperation and punishment in public goods experiments. *American Economic Review*. **90** 980-994.
- E. Fehr, Schmidt, K. 1999. A theory of fairness, competition and cooperation. *Quarterly Journal of Economics*. **114**(3) 817-868.
- U. Fischbacher. 2007. z-Tree: Zurich toolbox for ready-made economic experiments. *Experimental Economics*. **10**(2) 171-178.
- R. Fisman, Khanna, T. 1999. Is trust a historical residue? Information flows and trust levels. *Journal of Economic Behavior & Organization*. **38**(1) 79-92.
- G.R. Fréchette, Yuksel, S. 2013. *Infinitely Repeated Games in the Laboratory: Four Perspectives on Discounting and Random Termination*. Working paper.
- S. Frederick. 2005. Cognitive Reflection and Decision Making. *The Journal of Economic Perspectives*. **19**(4) 25-42.
- D. Fudenberg, Maskin, E.S. 1986. The Folk Theorem in Repeated Games with Discounting or with Incomplete Information. *Econometrica*. **54**(3) 533-554.
- D. Fudenberg, Maskin, E.S. 1990. Evolution and cooperation in noisy repeated games. *American Economic Review*. **80**(2) 274-279.
- D. Fudenberg, Peysakhovich, A. 2014. *Recency, records and recaps: learning and non-equilibrium behavior in a simple decision problem*. ACM, Palo Alto, California, USA.
- D. Fudenberg, Rand, D.G., Dreber, A. 2012. Slow to Anger and Fast to Forgive: Cooperation in an Uncertain World. *American Economic Review*. **102**(2) 720-749.
- S. Gächter, Herrmann, B., Thöni, C. 2004. Trust, voluntary cooperation, and socio-economic background: survey and experimental evidence. *Journal of Economic Behavior & Organization*. **55**(4) 505-531.
- S. Gächter, Herrmann, B., Thöni, C. 2010. Culture and cooperation. *Philosophical Transactions of the Royal Society B: Biological Sciences*. **365**(1553) 2651-2661.
- A.S. Gerber, Rogers, T. 2009. Descriptive Social Norms and Motivation to Vote: Everybody's Voting and so Should You. *The Journal of Politics*. **71**(1) 178-191.
- G. Gigerenzer, Goldstein, D.G. 1996. Reasoning the fast and frugal way: models of bounded rationality. *Psychological review*. **103**(4) 650.
- G. Gigerenzer, Todd, P.M., Group, A.R. 1999. *Simple heuristics that make us smart*. Oxford University Press, Oxford, UK.
- T. Gilovich, Griffin, D., Kahneman, D. 2002. *Heuristics and biases: The psychology of intuitive judgment*. Cambridge University Press.
- H. Gintis. 2003. The hitchhiker's guide to altruism: Gene-culture coevolution, and the internalization of norms. *Journal of theoretical biology*. **220**(4) 407-418.
- E.L. Glaeser, Laibson, D.I., Scheinkman, J.A., Soutter, C.L. 2000. Measuring Trust. *The Quarterly Journal of Economics*. **115**(3) 811-846.
- U. Gneezy, List, J.A. 2006. Putting behavioral economics to work: Testing for gift exchange in labor markets using field experiments. *Econometrica*. **74**(5) 1365-1384.
- K.E. Hauge, Brekke, K.A., Johansson, L.-O., Johansson-Stenman, O., Svedsäter, H. 2014. *Keeping others in our mind or in our heart? Distribution games under cognitive load*.
- J. Henrich, Boyd, R., Bowles, S., Camerer, C., Fehr, E., Gintis, H., McElreath, R., Alvard, M., Barr, A., Ensminger, J., Henrich, N.S., Hill, K., Gil-White, F., Gurven, M., Marlowe, F.W., Patton, J.Q., Tracer, D. 2005. "Economic man" in cross-cultural perspective: Behavioral experiments in 15 small-scale societies. *Behavioral and brain science*. **28** 795-855.

- J. Henrich, Ensminger, J., McElreath, R., Barr, A., Barrett, C., Bolyanatz, A., Cardenas, J.C., Gurven, M., Gwako, E., Henrich, N., Lesorogol, C., Marlowe, F., Tracer, D., Ziker, J. 2010. Markets, Religion, Community Size, and the Evolution of Fairness and Punishment. *Science*. **327**(5972) 1480-1484.
- J. Henrich, McElreath, R., Barr, A., Ensminger, J., Barrett, C., Bolyanatz, A., Cardenas, J.C., Gurven, M., Gwako, E., Henrich, N., Lesorogol, C., Marlowe, F., Tracer, D., Ziker, J. 2006. Costly Punishment Across Human Societies. *Science*. **312**(5781) 1767-1770.
- B. Herrmann, Thoni, C., Gächter, S. 2008. Antisocial punishment across societies. *Science*. **319**(5868) 1362-1367.
- H. Herz, Taubinsky, D. 2013. Market Experience is a Reference Point in Judgments of Fairness. Available at SSRN: <http://ssrn.com/abstract=2297773>
- J.J. Horton, Rand, D.G., Zeckhauser, R.J. 2011. The Online Laboratory: Conducting Experiments in a Real Labor Market. *Experimental Economics*. **14**(3) 399-425.
- J.J. Jordan, McAuliffe, K., Rand, D.G. 2014. Third-Party Punishment is Motivated by Anger and is not an Artifact of Self-Focused Envy or the Strategy Method. Available at SSRN: <http://ssrn.com/abstract=2427274>.
- J.J. Jordan, Peysakhovich, A., Rand, D.G. In press. *Why we cooperate*. MIT Press, Cambridge, MA.
- D. Kahneman. 2003. A perspective on judgment and choice: Mapping bounded rationality. *American Psychologist*. **58**(9) 697-720.
- R. La Porta, Lopez-de-Silanes, F., Shleifer, A., Vishny, R.W. 2001. Trust in large organizations. *Social Capital: A Multifaceted Perspective* 310.
- C.R. Leana, Buren, H.J.v., III. 1999. Organizational Social Capital and Employment Practices. *The Academy of Management Review*. **24**(3) 538-555.
- D.K. Levine. 1998. Modeling Altruism and Spitefulness in Experiments. *Review of Economic Dynamics*. **1**(3) 593-622.
- G.F. Loewenstein, O'Donoghue, T. 2004. Animal Spirits: Affective and Deliberative Processes in Economic Behavior. Available at SSRN: <http://ssrn.com/abstract=539843>.
- S. Lotz. 2014. Spontaneous Giving Under Structural Inequality: Intuition Promotes Cooperation in Asymmetric Social Dilemmas. Available at SSRN: <http://ssrn.com/abstract=2513498>.
- G.J. Mailath, Samuelson, L. 2006. Repeated games and reputations: long-run relationships. *OUP Catalogue*.
- D.J. McAllister. 1995. Affect-and cognition-based trust as foundations for interpersonal cooperation in organizations. *Academy of management journal*. **38**(1) 24-59.
- N. Nunn, Wantchekon, L. 2011. The slave trade and the origins of mistrust in Africa. *American Economic Review*. **101**(7) 3221-3253.
- E. Ostrom, Walker, J., Gardner, R. 1992. Covenants With and Without a Sword: Self-Governance is Possible. *The American Political Science Review*. **86**(2) 404-417.
- A. Ouss, Peysakhovich, A. 2013. When Punishment Doesn't Pay: 'Cold Glow' and Decisions to Punish. Available at SSRN: <http://ssrn.com/abstract=2247446>.
- E. Pedersen, Kurzban, R., McCullough, M. 2013. Do humans really punish altruistically? A closer look. *Proc R Soc B* <http://dx.doi.org/10.1098/rspb.2012.2723>.
- A. Peysakhovich, Nowak, M.A., Rand, D.G. 2014. Humans Display a 'Cooperative Phenotype' that is Domain General and Temporally Stable. *Nature Communications*. **5** 4939.
- A. Peysakhovich, Plagborg-Møller, M. 2012. A note on proper scoring rules and risk aversion. *Economics Letters*. **117**(1) 357-361.
- R.D. Putnam. 2000. *Bowling alone: The collapse and revival of American community*. Simon and Schuster.
- M. Rabin. 1993. Incorporating Fairness into Game Theory and Economics. *The American Economic Review*. **83**(5) 1281-1302.
- D.G. Rand, Fudenberg, D., Dreber, A. 2014a. It's the thought that counts: The role of intentions in noisy repeated games. Available at SSRN: <http://ssrn.com/abstract=2259407>.
- D.G. Rand, Greene, J.D., Nowak, M.A. 2012. Spontaneous giving and calculated greed. *Nature*. **489**(7416) 427-430.
- D.G. Rand, Newman, G.E., Wurzbacher, O. 2014b. Social context and the dynamics of cooperative choice. *Journal of Behavioral Decision Making* doi: 10.1002/bdm.1837.
- D.G. Rand, Nowak, M.A. 2011. The evolution of antisocial punishment in optional public goods games. *Nat Commun*. **2** 434.
- D.G. Rand, Nowak, M.A. 2012. Evolutionary dynamics in finite populations can explain the full range of cooperative behaviors observed in the centipede game. *Journal of theoretical biology*. **300** 212-221.
- D.G. Rand, Nowak, M.A. 2013. Human Cooperation. *Trends in Cognitive Sciences*. **17**(8) 413-425.
- D.G. Rand, Peysakhovich, A., Kraft-Todd, G.T., Newman, G.E., Wurzbacher, O., Nowak, M.A., Green, J.D. 2014c. Social Heuristics Shape Intuitive Cooperation. *Nature Communications*. **5** 3677.

- D.G. Rand, Tarnita, C.E., Ohtsuki, H., Nowak, M.A. 2013. Evolution of fairness in the one-shot anonymous Ultimatum Game. *Proceedings of the National Academy of Sciences*. **110**(7) 2581-2586.
- D.G. Rand, Yoeli, E., Hoffman, M. 2014d. Harnessing Reciprocity to Promote Cooperation and the Provisioning of Public Goods. *Policy Insights from the Behavioral and Brain Sciences*. **1**(1) 263-269.
- P.J. Richerson, Boyd, R. 2005. *Not by genes alone: How culture transformed human evolution*. University of Chicago Press, Chicago.
- S.G. Roch, Lane, J.A.S., Samuelson, C.D., Allison, S.T., Dent, J.L. 2000. Cognitive Load and the Equality Heuristic: A Two-Stage Model of Resource Overconsumption in Small Groups. *Organizational Behavior and Human Decision Processes*. **83**(2) 185-212.
- A.E. Roth. 2002. The economist as engineer: Game theory, experimentation, and computation as tools for design economics. *Econometrica*. **70**(4) 1341-1378.
- D.M. Rousseau, Sitkin, S.B., Burt, R.S., Camerer, C. 1998. Not so different after all: A cross-discipline view of trust. *Academy of management review*. **23**(3) 393-404.
- L. Samuelson. 2001. Introduction to the Evolution of Preferences. *Journal of Economic Theory*. **97**(2) 225-230.
- P. Sapienza, Zingales, L., Guiso, L. 2006. *Does culture affect economic outcomes?* National Bureau of Economic Research.
- J.F. Schulz, Fischbacher, U., Thöni, C., Utikal, V. 2014. Affect and fairness: Dictator games under cognitive load. *Journal of Economic Psychology*. **41** 77-87.
- A. Shenhav, Rand, D.G., Greene, J.D. 2012. Divine intuition: Cognitive style influences belief in God. *Journal of Experimental Psychology: General*. **141**(3) 423-428.
- J.P. Simmons, Nelson, L.D., Simonsohn, U. 2013. Life after P-Hacking. Available at SSRN: <http://ssrn.com/abstract=2205186>.
- G. Tabellini. 2007. The scope of cooperation: Values and incentives. *Innocenzo Gasparini Institute for Economic Research Working Paper*(328).
- R.H. Thaler, Benartzi, S. 2004. Save More Tomorrow™: Using behavioral economics to increase employee saving. *Journal of political Economy*. **112**(S1) S164-S187.
- R.H. Thaler, Sunstein, C.R. 2008. *Nudge: Improving decisions about health, wealth, and happiness*. Yale University Press.
- G. Tinghög, Andersson, D., Bonn, C., Böttiger, H., Josephson, C., Lundgren, G., Västfjäll, D., Kirchler, M., Johannesson, M. 2013. Intuition and cooperation reconsidered. *Nature*. **497**(7452) E1-E2.
- A. Tversky, Kahneman, D. 1974. Judgment under Uncertainty: Heuristics and Biases. *Science*. **185**(4157) 1124-1131.
- P.P.J.L. Verkoeijen, Bouwmeester, S. 2014. Does Intuition Cause Cooperation? *PLoS ONE*. **9**(5) e96654.
- M.B. Welsh, Burns, N.R., Delfabbro, P.H. 2013. *The Cognitive Reflection Test: how much more than Numerical Ability?* Cognitive Science Society.
- T. Yamagishi, Horita, Y., Mifune, N., Hashimoto, H., Li, Y., Shinada, M., Miura, A., Inukai, K., Takagishi, H., Simunovic, D. 2012. Rejection of unfair offers in the ultimatum game is no evidence of strong reciprocity. *Proceedings of the National Academy of Sciences*.
- E. Yoeli, Hoffman, M., Rand, D.G., Nowak, M.A. 2013. Powering up with indirect reciprocity in a large-scale field experiment. *Proceedings of the National Academy of Sciences*. **110**(Supplement 2) 10424-10429.

## **Habits of Virtue: Creating Norms of Cooperation and Defection in the Laboratory**

Alexander Peysakhovich and David G. Rand

### **Appendix**

1. Experiment 1: Stage A manipulation check .....	A2
2. Experiment 1: Principal component analysis of play in Stage B games .....	A3
3. Experiment 1: Effect of Stage A on Stage B prosociality .....	A5
4. Supplemental Experiment: Manipulating prosociality in the online laboratory.....	A10
4.1 Experimental design .....	A10
4.2 Results.....	A11
5. Experiment 2: Stage A manipulation check .....	A14
6. Experiment 2: Effect of Stage A on Stage B Third Party Punishment .....	A15
7. Experiment 2: Effect of Stage A on Stage B Public Goods Game Punishment .....	A20
8. Effect of Stage A on generalized trust .....	A21
9. Stage A Game Lengths.....	A23
10. Instructions: Experiment 1 and 2 Stage A (C-culture version) .....	A25
11. Instructions: Experiment 1 Stage B.....	A29
12. Instructions: Experiment 2 Stage B.....	A33
13. Instructions: Supplemental Experiment.....	A38

**1. Experiment 1: Stage A manipulation check**

**Table A1:** The effect of C-culture treatment on RPD cooperation in Stage A of Experiment 1. Shown are coefficients from OLS regression, with cooperation in the RPD as the dependent variable (0=D, 1=C; one observation per decision). Standard errors in parentheses clustered at session level and bootstrapped 1000x.

	PD Decision	PD Decision
1=C-Culture	0.554 (0.072)***	0.315 (0.049)***
Decision number		-0.001 (0.001)
C-Culture X Decision Num		0.009 (0.002)***
Constant	0.154 (0.038)***	0.179 (0.039)***
$R^2$	0.32	0.34
$N$	5,092	5,092

\*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$

Standard errors clustered at session level and bootstrapped 1000x.



## 2. Experiment 1: Principal component analysis of play in Stage B games

To provide support for our decision to jointly analyze the five transfer decisions as our measure of prosociality (and not the UG MAO), we use principal component analysis (PCA). Intuitively, PCA is a statistical technique used to reduce high dimensional data sets to a smaller number of dimensions in an optimal manner. The PCA analysis treats data points coming from  $n$  variables as vectors in an  $n$ -dimensional linear space and then computes a new set of basis vectors called “principal components” ordered by the amount of original variance retained if the data set were to be projected onto the principal component.

Applying PCA to our data can be thought of as asking whether the six decisions are really measuring independent aspects of behavior, or whether there are high correlations within subgroups of decisions. We perform a PCA using the standardized (correlation) matrix of the six decisions. The six resulting components have eigenvalues of 2.6, 1.23, 0.7, 0.5, 0.49 and 0.43. Consistent with our definition of prosociality, the first (most informative) principal component consists of a roughly equally weighted combination of all the transfer decisions: PGG, TG Trustor, TG Trustee, DG and UG Sender (Table A2); with almost no weight on the UG MAO. Thus, these results statistically validate our decision to treat the five transfer decisions as repeat measures of cooperation. Further support for this conclusion comes from Peysakhovich et al (2014), who find that an individual’s play in the DG, PGG and TG are strongly interrelated, but are not related to UG MAOs (importantly, these results, which were not preceded by any Stage A, show that the results of our PCA here are not caused by Stage A treatment effects).

**Table A2:** Loadings from the PCA on decisions in Experiment 1. To aid interpretation, loadings greater than 0.3 are bolded.

Component	1	2	3	4	5	6
PGG	<b>0.49</b>	-0.11	-0.20	0.07	<b>0.55</b>	<b>0.63</b>
DG	<b>0.47</b>	0.13	-0.27	<b>0.73</b>	-0.21	<b>-0.34</b>
TG Trustor	<b>0.46</b>	-0.22	<b>0.42</b>	-0.26	<b>0.39</b>	<b>-0.58</b>
TG Trustee	<b>0.47</b>	-0.11	<b>0.44</b>	-0.14	<b>-0.66</b>	<b>0.35</b>
UG Offer	<b>0.32</b>	<b>0.56</b>	<b>-0.47</b>	<b>-0.58</b>	-0.12	-0.13
UG MAO	-0.04	<b>0.77</b>	<b>0.54</b>	0.21	0.22	0.12
Variance explained	43%	21%	12%	9%	8%	7%

For completeness, we also show the pairwise correlations between the 6 game decisions (correlations which are significant at the  $p < 0.05$  level after Bonferonni correction for multiple comparisons are indicated by \*).

	PGG	DG	TG Trustor	TG Trustee	UG Offer	UG MAO
PGG	1					
DG	0.50*	1				
TG Trustor	0.50*	0.39*	1			
TG Trustee	0.46*	0.43*	0.53*	1		
UG Offer	0.31*	0.37*	0.18	0.23	1	
UG MAO	-0.13	0.01	-0.12	-0.06	0.23	1

The results paint a similar picture to the PCA: PGG, DG, and the two TG decisions are highly inter-correlated, UG Offer is correlated with these decisions but somewhat less strongly, and UG MAO is not correlated at all with the giving decisions.

When considering the relationship between these different games, it is also interesting to note that certain transfer decisions have both a self-interested component and a prosocial component (TG Trustor and UG Offer), while others have only a prosocial component (PGG, DG, TG Trustee). That is to say that a self-interested player might transfer as the TG Trustee or make a non-zero UG Offer, depending on their expectations about the behavior of their co-player. Self-interested players would never transfer anything, however, in the PGG, DG and TG Trustee decisions. This distinction is interesting because we find that Stage A has a larger effect on the purely prosocial decisions than on the decisions that have a self-interest component (see main text Figure 2). Thus this suggests that Stage A is particularly affecting subjects' level of prosociality, rather than their self-interested calculations.

### 3. Experiment 1: Effect of Stage A on Stage B prosociality

**Table A3:** Regression analyses (OLS) of effect of Stage A treatment on Stage B prosociality. Dependent variable is fraction transferred in Stage B decision, with five observations per subject (one for each of the five Stage B transfer decisions). When including decision-type dummies, DG is held out as the baseline.

	(1)	(2)	(3)
1=C-Culture	0.100 (0.026)***	0.100 (0.028)***	0.132 (0.041)**
1=PGG		0.167 (0.022)***	0.167 (0.007)***
1=TG Trustor		0.344 (0.031)***	0.366 (0.032)***
1=TG Trustee		0.119 (0.018)***	0.126 (0.023)***
1=UG Proposer		0.209 (0.025)***	0.255 (0.016)***
C-Culture X PGG			0.001 (0.050)
C-Culture x TG Trustor			-0.048 (0.066)
C-Culture X TG Trustee			-0.016 (0.039)
C-Culture X UG Proposer			-0.099 (0.041)*
Constant	0.278 (0.007)***	0.110 (0.011)***	0.095 (0.007)***
$R^2$	0.02	0.13	0.13
$N$	480	480	480

\*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$

Standard errors clustered at session level and bootstrapped 1000x

**Table A4.** Regression analyses (OLS) of effect of Stage A treatment on Stage B behavior by game. Dependent variable is fraction transferred in the indicated game (or, the final two columns, predicted fraction returned for TG Trustor Beliefs and fraction demanded for UG MAO).

	DG	PGG	TG Trustor	TG Trustee	UG Proposer	TG Trustor Beliefs	UG MAO
C-Culture	0.132 (0.041)**	0.133 (0.035)***	0.084 (0.098)	0.116 (0.020)***	0.033 (0.016)*	0.153 (0.042)***	-.041 (.037)
Constant	0.095 (0.007)***	0.262 (0.001)***	0.462 (0.040)***	0.221 (0.016)***	0.350 (0.014)***	0.334 (0.011)***	.256 (0.009)***
$R^2$	0.09	0.04	0.01	0.05	0.01	0.08	0.02
$N$	96	96	96	96	96	96	96

\*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$

Standard errors clustered at session level and bootstrapped 1000x.

**Table A5:** Effect of Stage A treatment on Stage B prosociality in Experiment 1, interacted with (and separately analyzed by) number of correct answers on the Cognitive Reflection Test (CRT). Dependent variable is fraction transferred in Stage B decision, with five observations per subject (one for each of the five Stage B transfer decisions).

	All subjects	Intuitive subjects (>0 Incorrect responses)	Reflective subjects (0 Incorrect responses)
1=C-Culture	0.227 (0.077)**	0.160 (0.062)**	-0.029 (0.035)
# Correct CRT answers	0.057 (0.029)*		
C-Culture X CRT	-0.070 (0.032)*		
1=DG	-0.167 (0.022)***	-0.155 (0.025)***	-0.194 (0.034)***
1=TG Trustor	0.177 (0.051)***	0.142 (0.069)*	0.250 (0.057)***
1=TG Trustee	-0.049 (0.020)*	-0.038 (0.013)**	-0.072 (0.068)
1=UG Proposers	0.042 (0.023)	0.035 (0.031)	0.058 (0.033)
Constant	0.175 (0.057)**	0.253 (0.019)***	0.347 (0.032)***
$R^2$	0.15	0.14	0.18
$N$	480	325	155

\*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$

Standard errors clustered at session level and bootstrapped 1000x.

**Table A6:** Effect of Stage A treatment on Stage B prosociality in Experiment 1, interacted with (and separately analyzed by) number of intuitive answers on the Cognitive Reflection Test (CRT). Dependent variable is fraction transferred in Stage B decision, with five observations per subject (one for each of the five Stage B transfer decisions).

	All subjects	Intuitive subjects (>0 Intuitive responses)	Reflective subjects (0 Intuitive responses)
1=C-Culture	0.032 (0.020)	0.171 (0.064)**	-0.003 (0.031)
# Intuitive CRT answers	-0.064 (0.029)*		
C-Culture X Intuitive CRT	0.078 (0.039)*		
1=DG	-0.167 (0.022)***	-0.153 (0.029)***	-0.189 (0.035)***
1=TG Trustor	0.177 (0.052)***	0.146 (0.069)*	0.223 (0.053)***
1=TG Trustee	-0.049 (0.020)*	-0.051 (0.021)*	-0.046 (0.061)
1=UG Proposers	0.042 (0.023)	0.027 (0.041)	0.064 (0.031)*
Constant	0.331 (0.029)***	0.253 (0.035)***	0.317 (0.034)***
$R^2$	0.14	0.14	0.16
$N$	480	285	195

\*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$

Standard errors clustered at session level and bootstrapped 1000x.

**Table A7:** Regression analyses (OLS) of effect of Stage A treatment on Stage B prosociality, restricting to subjects from the three sessions of Experiment 1 in which mood was measured. Dependent variable is fraction transferred in Stage B decision, with five observations per subject (one for each of the five Stage B transfer decisions). When including decision-type dummies, DG is held out as the baseline. Standard errors are clustered at the level of the subject instead of the session, because only the mood measure was only administered in three sessions. Clustering at the session level does not qualitatively change the results.

	(1)	(2)	(3)	(4)
1=C-Culture	0.128 (0.066)	0.139 (0.068)*	0.128 (0.067)	0.139 (0.068)*
Mood (1 to 5)		-0.064 (0.045)		-0.064 (0.045)
1=PGG			-0.128 (0.044)**	-0.128 (0.044)**
1=TG Trustor			0.265 (0.057)***	0.265 (0.057)***
1=TG Trustee			-0.041 (0.038)	-0.041 (0.038)
1=UG Proposer			0.042 (0.050)	0.042 (0.050)
Constant	0.282 (0.043)***	0.494 (0.151)**	0.255 (0.054)***	0.466 (0.155)**
$R^2$	0.03	0.05	0.18	0.19
$N$	240	240	240	240

\*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$

Standard errors clustered at subject level

#### 4. *Supplemental Experiment: Manipulating prosociality in the online laboratory*

Our Supplemental Experiment examines the generalizability of our Experiment 1 results. To do so, we deploy an adapted version of our Stage A manipulation using the online labor market Amazon Mechanical Turk (Amir et al. (2012), Horton et al. (2011); MTurk) and replicate the Stage A effect on giving in a DG where the recipient does not participate in Stage A and has no interactions other than receiving money from the dictator.<sup>1</sup> This experiment provides further evidence that Stage A alters social preferences rather than expectations about the behavior of (or beliefs about the behavioral type of) Stage B interaction partners, and shows that our effect is not unique to undergraduate students.

##### 4.1 *Experimental design*

Due to technical limitations of the MTurk platform, it is difficult to have simultaneous play with live feedback on MTurk. Thus we have to adapt our Stage A design to MTurk's constraints as follows.

Subjects are recruited and told that they will play a series of RPD games (continuation probability 0.8) matched with other subjects from MTurk. Approximately 35% of subjects are put into a "strategy choice" condition where they choose a strategy to be used in all of their RPD games. Two possible strategies are offered: Tit-for-Tat or Always Defect. The remaining 65% of subjects are assigned to "active play" conditions. These subjects play 4 random length RPDs (as in the laboratory experiments, the same set of game lengths was used for all subjects) matched against the strategy choices of subjects in the "strategy choice" condition (and affecting the payoffs of those subjects). Thus, no deception is used in this experiment.

To replicate the cooperation-supporting conditions of the C-Culture treatment, "Active play" subjects in the C-Culture treatment are matched against four "strategy choice" players who had chosen Tit-for-Tat. Thus "active play" subjects in this condition can experience mutual cooperation, but are prevented from exploiting their opponents. To replicate the defection-supporting conditions of the D-Culture treatment, "Active play" subjects in the D-culture treatment, conversely, are matched with four "strategy choice" subjects who had chosen Always Defect, so no mutual cooperation can arise.

---

<sup>1</sup> The MTurk subject pool is substantially more diverse than the college undergraduates used in Experiments 1 through 3 on a number of dimensions including age, education and income. Although one might worry about the lack of control implied by online experiments, a number of studies have demonstrated the validity of economic game experiments conducted on MTurk (Amir et al. (2012), Horton et al. (2011)). MTurk subjects have also been shown to not differ from a nationally representative sample on various psychological measures (Paolacci et al. (2010), Simons & Chabris (2012)).



To control for comprehension, an important issue on MTurk, all subjects take a quiz on the rules of the RPD interaction. Subjects that answer incorrectly are not allowed to participate. We analyze the ex-post behavior of the 237 “active play” subjects who passed the RPD comprehension questions.

After completing the four RPD games, subjects play a single-shot DG. Commensurate with standard wages on MTurk (Rand et al. (2012), Amir et al. (2012), Horton et al. (2011)), subjects are endowed with 80 cents that they choose how to split between themselves and another individual. Prior to making their decision, subjects are asked to complete a comprehension question regarding the payoff structure of the DG; as shown below, our results are robust to controlling for comprehension or excluding the 25% of subjects that failed the comprehension check. To additionally control for attention, individuals are asked to re-indicate the decision they made on the next screen of the study. All participants did so correctly.

Importantly, subjects are informed that their recipient in the DG had not participated in the RPDs, and receives no earnings other than what is given to them by the subject in the DG. Thus, changes in DG behavior here reflect changes in social preferences targeted at complete strangers with whom subjects have had no previous interaction.

As the Supplemental Experiment uses a more heterogeneous population than in our experiments using college students, we provide basic descriptive statistics in Table A8.

**Table A8.** Descriptive statistics for the Supplemental Experiment.

Age	M=29.51 s.d. = 10.04
Gender	63% male
Completed High school	98.4%
Completed at least BA degree	77%
Median reported income in 2012	\$15,000-\$25,000

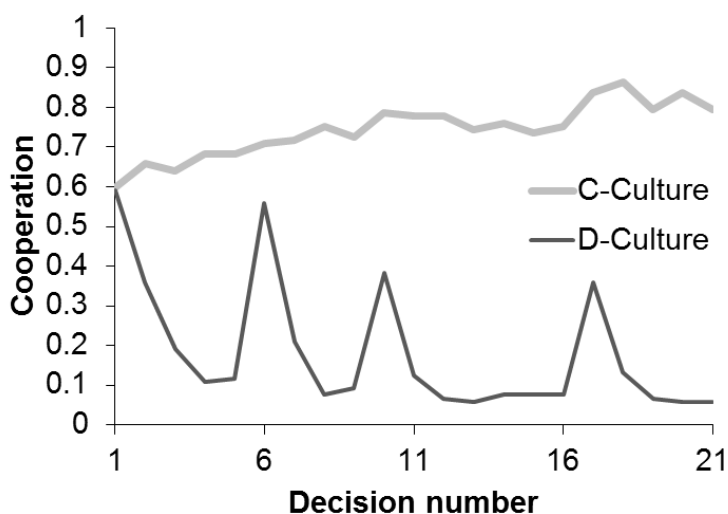
## 4.2 Results

We first show that our modified Stage A successfully creates environments of persistent cooperation and defection. Indeed, as can be seen in Figure A1, subjects learn to cooperate in the C-Culture treatment (i.e. when matched with four TFT partners) and learn to defect in the D-Culture treatment (i.e. when matched with four ALLD partners).<sup>2</sup> As a result, there is significantly more

---

<sup>2</sup> The cooperation spikes in the D-Culture time series reflect the decisions in which new games with new partner began: as explained above, the same sequence of game lengths was used across all sessions.

cooperation in our C-culture treatment than our D-culture, as well as the expected learning effects (Table A9).



**Figure A1:** Cooperation in the PD as a function of decision number.

**Table A9:** Regression analysis (OLS) of time series of cooperation in C-Culture and D-Culture treatments. Dependent variable is cooperation in the RPD (0=D, 1=C; one observation per decision).

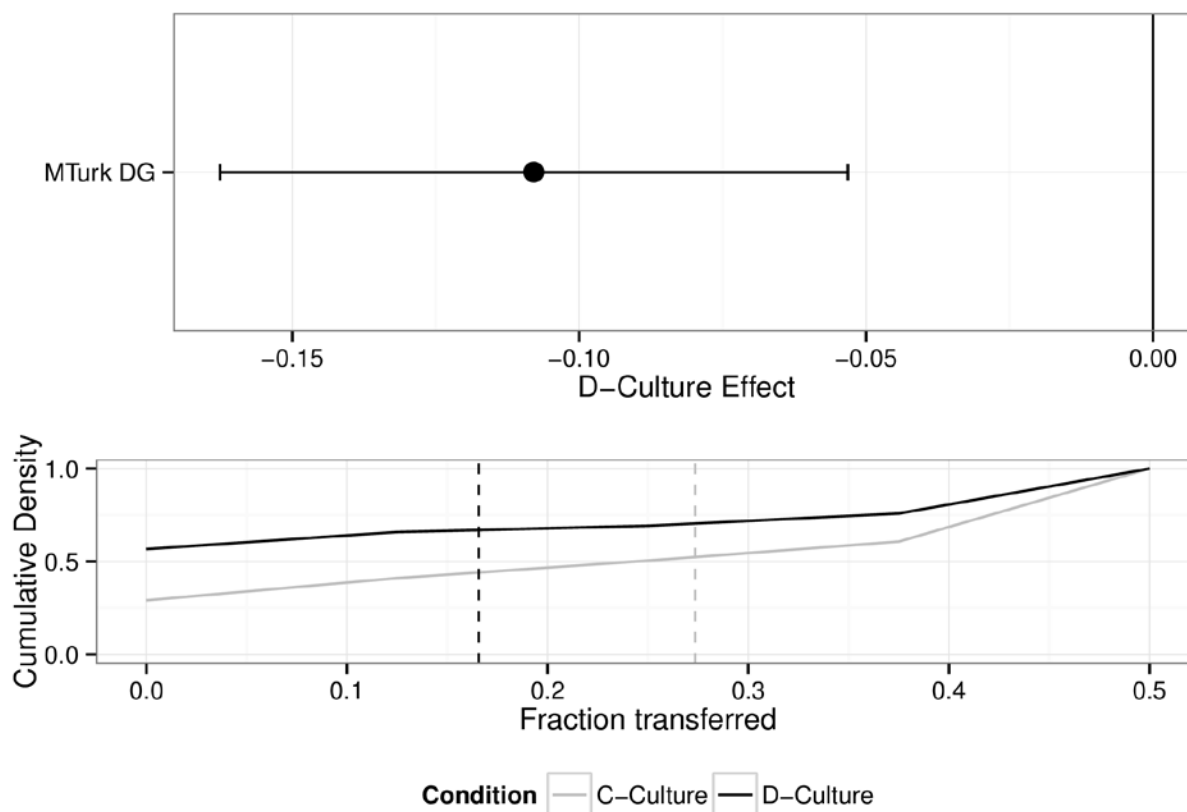
	PD Decision	PD Decision
1=C-Culture	0.562 (0.037)***	0.315 (0.048)***
Decision number		-0.011 (0.001)***
C-Culture X Decision Number		0.017 (0.002)***
Constant	0.183 (0.020)***	0.326 (0.027)***
$R^2$	0.32	0.34
$N$	4,977	4,977

\*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$

Standard errors in parentheses clustered at the subject level.

We now turn to examining play in the one-shot DG of Stage B. There is a strong effect of Stage A experience, with subjects in the D-Culture giving substantially less in the DG compared to subjects from the C-Culture (mean transfer in D-Culture 16% of endowment vs C-Culture 27% of endowment; Figure A2). Table A10 shows that there is a strong significant effect of treatment on DG behavior, both with and without controls, or exclusions for comprehension. Thus we successfully replicate the effect of immersion in a cooperative or non-cooperative environment on one-shot anonymous prosociality using

our online sample. This replication demonstrates that the effect is not unique to college undergraduates, and that behaviors induced by immersion in our C vs. D-Cultures are exhibited in decisions involving complete passive outsiders with whom subjects have never previously interacted.



**Figure A2:** Top panel shows effect of D-Culture on DG transfer (error bars indicate 95% confidence intervals). Bottom panel shows cumulative distribution function of DG transfers by treatment.

**Table A10:** Regression analysis (OLS) of the effects of Stage A treatment on DG transfer. Dependent variable is number of cents transferred (out of 80 cents).

	All subjects	All subjects	Excluding non-comprehenders
1=C-Culture	8.630 (2.233)***	8.669 (2.255)***	7.276 (2.630)**
Failed comprehension		-0.356 (2.593)	
Constant	13.250 (1.569)***	13.321 (1.656)***	13.958 (1.779)***
$R^2$	0.06	0.06	0.04
$N$	237	237	177

\*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$

Robust standard errors in parentheses.

### 5. Experiment 2: Stage A manipulation check

**Table A11:** The effect of Stage A parameters on PD cooperation in Stage A of Experiment 2. Shown are coefficients from OLS (linear probability model) regression with cooperation in the RPD as the dependent variable (0=D, 1=C; one observation per decision).

	PD Decision	PD Decision
1=C-Culture	0.358 (0.087)***	0.178 (0.101)
Decision number		-0.006 (0.003)*
C-Culture X Decision Num		0.010 (0.003)***
Constant	0.217 (0.058)***	0.321 (0.082)***
$R^2$	0.13	0.14
$N$	4,616	4,616

\*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$

Standard errors in parentheses clustered at session level and bootstrapped 1000x.

We note that there is less cooperation in the Stage A C-Culture in Experiment 2 than there was in the corresponding condition of Experiment 1. This is because in Experiment 2, two of the C-culture sessions contained a non-trivial fraction of persistent defectors who drove cooperation rates down for the whole group. Bootstrapping the data at the session level addresses this session-level variation.

### 6. Experiment 2: Effect of Stage A on Stage B Third Party Punishment

**Table A12:** Regression analyses (OLS) of effect of Stage A treatment on Stage B third party punishment in the 3PDG and 3PPD games. Dependent variable in cols 1-5 is the fraction of the punishment endowment spent, with 8 observations per subject in cols 1-3 (6 levels of punisher prosociality in the 3PDG and 2 levels of punisher prosociality in the 3PPD), 6 observations per subject in col 4, and 2 observations per subject in col 5. Dependent variable in col 6 is fraction of the punishment endowment spent on rewarding in the 3PPD game, with 2 observations per subject.

	Both games	Both games	Both games	3PDG	3PPD	3PPD Rewarding
I=C-Culture	1.940 (0.885)*	1.940 (0.892)*	1.935 (0.940)*	1.936 (0.967)*	1.953 (1.134)†	-0.139 (0.577)
Punisher's prosociality	-5.601 (0.816)***	-5.601 (0.861)***	-5.601 (0.858)***	-6.184 (1.068)***	-4.786 (1.247)***	4.607 (0.816)***
C-Culture X Punisher's prosociality	-2.947 (1.088)**	-2.947 (1.138)**	-2.947 (1.141)**	-2.951 (1.371)*	-2.942 (1.433)*	-0.850 (1.233)
Game		-0.393 (0.337)	-0.405 (0.230)			
C-Culture X Game			0.021 (0.679)			
Constant	7.122 (0.475)***	7.220 (0.463)***	7.223 (0.414)***	7.514 (0.161)***	6.411 (0.965)***	1.518 (0.429)***
$R^2$	0.18	0.18	0.18	0.16	0.22	0.12
$N$	976	976	976	732	244	244

† $p < 0.10$ ; \*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$

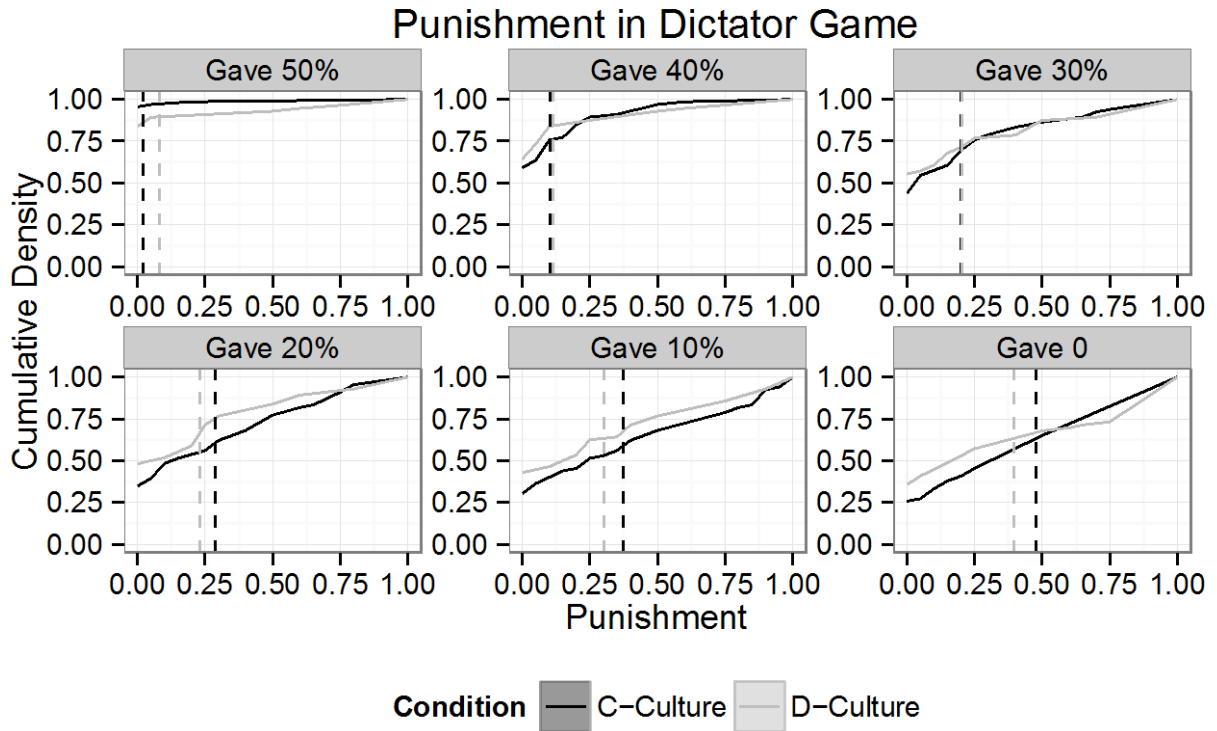
Standard errors clustered at session level and bootstrapped 1000x.

**Table A13:** Regression analyses (OLS) of effect of Stage A treatment on Stage B ratings of inappropriateness of punisher's action in the 3PDG and 3PPD games. Dependent variable is level of inappropriateness (using a 1-to-7 scale), with 8 observations per subject (6 levels of punisher prosociality in the 3PDG and 2 levels of punisher prosociality in the 3PPD).

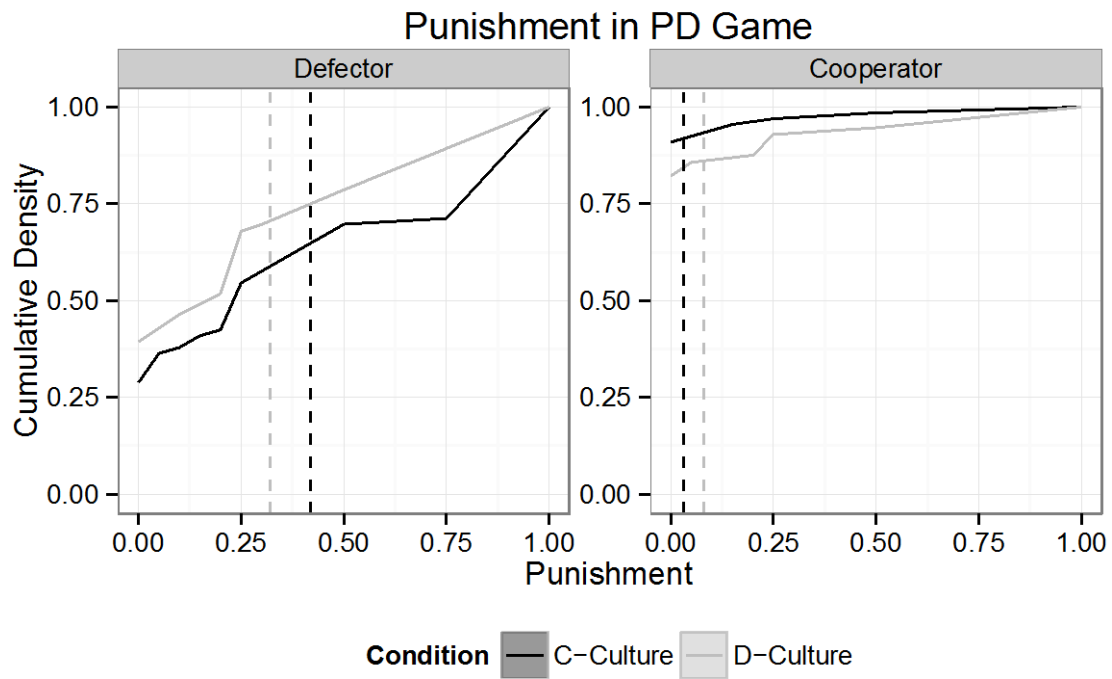
	Both games	Both games	Both games	3PDG	3PPD
1=C-Culture	0.486 (0.346)	0.486 (0.316)	0.464 (0.318)	0.458 (0.254)	0.560 (0.579)
Punisher's prosociality	-2.295 (0.229)***	-2.295 (0.220)***	-2.295 (0.228)***	-2.531 (0.150)***	-1.964 (0.603)**
C-Culture X Punisher's prosociality	-0.444 (0.474)	-0.444 (0.441)	-0.444 (0.468)	-0.433 (0.383)	-0.460 (0.838)
Game		-0.048 (0.110)	-0.095 (0.209)		
C-Culture X Game			0.088 (0.247)		
Constant	4.344 (0.194)***	4.356 (0.173)***	4.368 (0.155)***	4.486 (0.048)***	4.107 (0.493)***
$R^2$	0.23	0.23	0.23	0.22	0.25
$N$	976	976	976	732	244

† $p < 0.10$ ; \*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$

Standard errors clustered at session level and bootstrapped 1000x.



**Figure A3.** Cumulative density function (CDF) plots of fraction of punishment endowment spent in the 3PDG on punishment of each level of dictator transfer (from giving 50% to giving nothing), by treatment. Mean punishment amount shown with vertical dotted lines. We see that people from the C-Culture punish selfish dictators more, and punish fair dictators less.



**Figure A4.** CDF plots of fraction of punishment endowment spent in the 3PPD on punishment of cooperators and defectors, by treatment. Mean punishment amount shown with vertical dotted lines. We see that people from the C-Culture punish defectors more, and punish cooperators less.



**Table A14:** Regression analyses (OLS) of effect of Stage A treatment on Stage B third party punishment in the 3PDG and 3PPD games, by responses to post-experiments Cognitive Reflection Test (CRT). Dependent variable is the fraction of the punishment endowment spent, with 8 observations per subject (6 levels of punishee prosociality in the 3PDG and 2 levels of punishee prosociality in the 3PPD).

Classification method	# of Incorrect CRT responses		# of Intuitive CRT responses	
Subject type	Intuitive (>0 Incorrect)	Reflective (0 Incorrect)	Intuitive (>0 Intuitive)	Reflective (0 Intuitive)
1=C-Culture	2.034 (1.078)	0.031 (2.429)	2.370 (1.267)	1.181 (2.160)
Punishee's prosociality	-4.688 (0.909)***	-7.696 (1.149)***	-4.716 (1.125)***	-6.970 (1.194)***
C-Culture X Punishee's prosociality	-2.961 (1.470)*	-1.569 (2.262)	-3.603 (1.615)*	-2.070 (1.993)
Constant	6.472 (0.129)***	8.613 (1.429)***	6.850 (0.403)***	7.542 (1.544)***
$R^2$	0.14	0.22	0.16	0.22
$N$	648	224	632	344

\*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$

Standard errors clustered at session level and bootstrapped 1000x.

### 7. Experiment 2: Effect of Stage A on Stage B Public Goods Game Punishment

**Table A15:** Regression analysis (OLS) of punishment decisions in PGP games. Dependent variable is the fraction of the punishment endowment spent, with 9 observations per subject total (9 different contribution levels of the punishee); how those 9 observations are split between the PSP and ASP regressions depends on the punisher's own contribution level. Contribution difference is normalized to range from -1 (punisher contributed nothing and punishee contributed everything) to 1 (punisher contributed everything and punishee contributed nothing).

	PSP		ASP	
	[Contribution Diff] > 0		[Contribution Diff] ≤ 0	
	(1)	(2)	(3)	(4)
1=C-Culture	0.491 (1.354)	-0.346 (1.202)	-1.150 (0.814)	-0.293 (0.406)
Contribution Diff.	11.379 (1.746)***	10.501 (2.987)***	-0.435 (0.904)	-1.781 (1.663)
C-Cult X Cont Diff		1.631 (3.520)		2.764 (1.672)†
Constant	2.140 (0.902)*	2.598 (0.917)**	1.470 (0.633)*	1.028 (0.352)**
$R^2$	0.17	0.17	0.03	0.04
$N$	271	271	339	339

†  $p < 0.10$ ; \*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$

Standard errors clustered at session level and bootstrapped 1000x.

### 8. Effect of Stage A on generalized trust

**Table A16:** Regression analysis (OLS) of responses to the generalized trust question adapted from the World Values Survey. Dependent variable is agreement with the statement “Most people can be trusted” (1-to-5 scale).

	All subjects	All subjects	Lab subjects	MTurk subjects
1=C-Culture	0.413 (0.119)***	0.412 (0.108)***	0.416 (0.202)*	0.409 (0.127)**
1=MTurk		0.259 (0.119)*		
Constant	2.845 (0.091)***	2.690 (0.125)***	2.687 (0.175)***	2.950 (0.097)***
$R^2$	0.04	0.05	0.04	0.04
$N$	394	394	157	237

\*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$

Standard errors clustered at session level and bootstrapped 1000x. As MTurk subjects are each independent observations, each subject from the MTurk experiment is treated as its own session.

**Table A17:** Regression analyses (OLS) of effect of Stage A on responses to the generalized trust question adapted from the World Values Survey, by responses to post-experiments Cognitive Reflection Test (CRT). Dependent variable is agreement with the statement “Most people can be trusted” (1-to-5 scale). Includes all subjects from Experiment 1 and Experiment 2 who were administered the WVS trust measure, but not subjects from the Supplemental Experiment (as the CRT was not included in the Supplemental Experiment).

Classification method	# of Incorrect CRT responses		# of Intuitive CRT responses	
Subject type	Intuitive (>0 Incorrect)	Reflective (0 Incorrect)	Intuitive (>0 Intuitive)	Reflective (0 Intuitive)
1=C-Culture	0.451 (0.214)*	0.325 (0.314)	0.585 (0.236)*	0.147 (0.295)
Constant	2.638 (0.155)***	2.818 (0.253)***	2.615 (0.159)***	2.821 (0.227)***
$R^2$	0.04	0.02	0.07	0.00
$N$	114	43	97	60

\*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$

Standard errors clustered at session level and bootstrapped 1000x.

We note that the moderation is strong when classifying subjects based on the number of specifically intuitive responses they gave (cols 3 vs 4), with nearly 4 times as large an effect among intuitive subjects compared to reflective subjects; but much weaker when classifying subjects based on the number of incorrect CRT responses (cols 1 and 2). This is perhaps not so surprising, given that the classification in cols 3 and 4 is a more direct measure of heuristic use than that in cols 1 and 2.

### 9. Stage A Game Lengths

Game lengths for each interaction were generated using the indicated continuation probability for each condition.

**Table A18.** Number of rounds per interaction of the C-Culture and D-Culture in Experiment 1.

Interaction #	C-Culture $\delta=7/8$	D-Culture $\delta=1/8$
1	12	1
2	1	2
3	3	1
4	2	1
5	1	1
6	9	1
7	6	1
8	7	2
9	5	2
10	6	2
11		1
12		1
13		2
14		1
15		1
16		1
17		1
18		1
19		1
20		1
21		1
22		1
23		1
24		1
25		1
26		1
27		1
28		1
29		1
30		1
31		2
32		2
33		1
34		1
35		2
36		1
37		1
38		1
39		1
40		1
41		1
42		2
43		1
44		1
45		1

**Table A19.** Number of rounds per interaction of the C-Culture and D-Culture in Experiment 1.

Interaction #	C-Culture $\delta=7/8$	D-Culture $\delta=1/8$
1	12	1
2	1	2
3	3	1
4	2	1
5	1	1
6	9	1
7	6	1
8		2
9		2
10		2
11		1
12		1
13		2
14		1
15		1
16		1
17		1
18		1
19		1
20		1
21		1
22		1
23		1
24		1
25		1
26		1
27		1
28		1
29		1

### ***10. Instructions: Experiment 1 and 2 Stage A (C-culture version)***

*Please read the following instructions carefully. If you have any questions, do not hesitate to ask us. Aside from this, no communication is allowed during the experiment.*

#### **Instructions**

This is a computerized experiment on decision-making. You will be paid for participating and the amount you earn will depend on the decisions that you make.

The full experiment should about 60 minutes. At the end of the experiment, you will be paid privately and in cash for your participation.

All information collected in this experiment will be anonymous and neither the experimenter nor other subjects will be able to link your identity to your decisions. In order to maintain this privacy, please do not reveal your decisions to any other subject.

We consider ourselves bound by the promises we are making to you in this protocol, we will do everything we say and there will be no surprises or tricks. We are interested in individual choices so please **remember that there are no right or wrong answers**.

#### **Payment**

In this experiment you will earn Monetary Units (MUs) through the decisions that you make. At the end of the experiment, these MUs will be converted into dollars at a rate of 30 MU per dollar.

In addition to any money you earn from your decisions you will also receive a \$10 show-up fee.

## The Experiment

The experiment will be split into two parts, Part A and Part B. *Your decisions in Part A will not at all affect what will happen in Part B.* You will be able to earn MU in both parts and your final income will be the sum of the MU you get in Part A and Part B. The instructions for Part B will be given on your computer screen after the completion of Part A.

### Part A

Part A will consist of many individual sub-sections called *interactions*.

At the beginning of each *interaction* you will be matched with another individual in the room. You will then play a random number of *rounds* with this individual.

In each *round*, you will both be presented with two options: **A** and **B**. You will both receive an income in MU at the end of each round. This income will depend on both the choices that you and your opponent have made.

The figures below show how your incomes in a round will be determined:

**FIGURE 1: YOUR INCOMES**

	<b>Other Player Chooses A</b>	<b>Other Player Chooses B</b>
<b>You choose A</b>	You get 4	You get 0
<b>You choose B</b>	You get 5	You get 1

**FIGURE 2: OTHER PLAYER'S INCOMES**

	<b>Other Player Chooses A</b>	<b>Other Player Chooses B</b>
<b>You choose A</b>	Other player gets 4	Other player gets 5
<b>You choose B</b>	Other player gets 0	Other player gets 1

*For example:*

If you choose A and the other person chooses A, you would both get 4 units.

If you choose A and the other person chooses B, you would get 0, and they would get 5 units.

If you choose B and the other person chooses A, you would get 5 units, and they would get 0 units.

If you choose B and the other person chooses B, you would both get 1 unit.



*You must make your choice within 30 seconds, if you do not, the computer will make a random choice for you for that round.*

Your income for each round will be calculated and presented to you on your computer screen. This income will be added to your current stock of MU.

Remember that the total number of MU you have at the end of the session will determine how much money you earn, at an exchange rate of 30 units = \$1.

### **Rounds Per Interaction**

After each *round*, there is a  $7/8$  probability of another *round*, and  $1/8$  probability that the *interaction* will end. This probability does not depend on how many rounds you have already played.

Once the interaction ends, you will be randomly re-matched with a different person in the room for another interaction. Choices that you make will not influence either the number of interactions you have or the number of rounds in any interaction.

### **Summary of Part A**

To summarize: every *interaction* you have with another person in the experiment includes a random number of rounds. In each *round* both you and the person you are matched with will make one decision, this decision will determine your incomes for that round.

In every interaction, after every *round*, there is a  $7/8$  probability of another *round*. There will be a number of interactions, and your behavior has no effect on the number of rounds or the number of interactions.

At the beginning of the experiment, you will start with some MU in your account. You will earn MU in every round of every interaction. At the end of the experiment, you will receive \$1 for every 30 MU in your account.

You will now take a very short quiz to make sure you understand the setup.

Part A will begin with one practice round. This round will not count towards your final income.

### QUIZ

If you choose <b>A</b> and the other person chooses <b>B</b> you will receive...	<ul style="list-style-type: none"> <li>a. 0 units</li> <li>b. 4 units</li> <li>c. 1 unit</li> </ul>
If you choose <b>B</b> and the other person choose <b>A</b> ,you will receive...	<ul style="list-style-type: none"> <li>a. 5 units</li> <li>b. 1 unit</li> <li>c. 0 units</li> </ul>
The number of rounds in an interaction depends on your actions in that interaction or other interactions.	<p style="text-align: center;">TRUE FALSE</p>
If you have already played 2 rounds, the probability that there will be another round in your interaction is...	<ul style="list-style-type: none"> <li>a. 4/8</li> <li>b. 7/8</li> <li>c. 0</li> <li>d. 1/8</li> </ul>
If you have already played 5 rounds, the probability that there will be another round in your interaction is...	<ul style="list-style-type: none"> <li>a. 4/8</li> <li>b. 7/8</li> <li>c. 0</li> <li>d. 1/8</li> </ul>

*If you have any questions, please ask right now. Aside from this, no communication is allowed during the experiment.*

## ***11. Instructions: Experiment 1 Stage B***

### **Screen 1 – General Instructions**

You will now be asked to make decisions in 4 different interactions. Unlike the interactions you have experienced so far each of these interactions will last for exactly one round.

For each of these interactions, you will be matched anonymously with new individuals that are currently in the experiment. The choices you make in each interaction will not affect subsequent interactions in any way.

Some interactions will include multiple roles (Player A and Player B), you will be asked to make choices for each role.

After you finish all the interactions, one interaction will be chosen randomly to actually be played. If this interaction had different roles, they will be randomly assigned. Only the outcome of the chosen interaction will affect your (and the other individuals') earned MUs. You will be shown your income from the randomly chosen interaction.

Because you do not know which of your decisions will be chosen it is in your interest to treat each decision as if it counts for real MU.

### **Screen 2 - PGG Instructions**

You are matched with a group of 3 other individuals. Each person in your group is endowed with 100 units. Each individual decides how many of the 100 units they are going to contribute to a common project that benefits all group members (as described below), and how many of them to keep for themselves.

The contributions of all 4 players are added up. The total sum is multiplied by 1.6 and then evenly split among all 4 players. Each player gets the same share from the project.

In addition to your earnings from the project, you also receive the units you chose not to contribute.

Thus, your income is:

$$100 - (\text{your contribution to the project}) + 1.6 \times (\text{sum of all contributions}) / 4$$

Example 1:

Suppose each person contributes 100 MU to the project. This 400 MU is multiplied by 1.6 and then split between all 4 people.

This means each person will receive 160 MU.

Example 2:

Suppose 3 people contribute 100 MU and 1 person contributes 0 MU. The 3 people that contributed will get 120 MU

The person that did not contribute will get  
 $120$  (from project) +  $100$  (that they kept) =  $220$  MU

### Screen 3 – PGG Contribution

You start with  $100$  MU.

How much will you contribute to the project? [Textbox]

### Screen 4 – TG Instructions

You will be matched with one other individual. One of you will be assigned the role of Player A, the other one of Player B.

In this interaction both players will start with  $50$  MUs. Player A will then be given a choice to send either  $0$  or  $50$  MU to Player B. If player A chooses to send, this MU will be tripled and Player B will receive  $150$  MU. Player B will then get to choose how many MU out of that  $150$  to send back to Player A.

At the end of the interaction, Player A's income will be:  $50 - (\text{How much Player A sent}) + (\text{How much Player B sent back})$

Player B's income will be:  $50 + (3 * \text{How much Player A sent}) - (\text{How much Player B sent back})$

Example

Player A chooses to send

Player B receives  $150$  units

Player B chooses to send back  $30$  units

Then Player A receives:  $50 - 50 + 30 = 30$  units

Player B receives:  $50 + 150 - 30 = 170$  units

### Screen 5 – TG Player A

You start with  $50$  MU.

If you are player A, would you like to send your MU to player B?

[Button = Send] [Button=Don't Send]

### Screen 5 – TG Player B

You start with  $50$  MU.

If player A chooses to send  $50$  MU (which will be tripled to  $150$ ), how much would you like to send to player A (up to  $150$ )?

[Textbox]

### Screen 6 – TG Beliefs

This question will not count for any of your payoffs, nor will its answer be revealed to any other subjects.

If you choose to send 50 MU, how many MU do you think other individuals in the experiment would send back on average (between 0 and 150)?

[Textbox]

### **Screen 7 – DG Instructions**

You will be matched with one other individual. One of you will be assigned the role of Player A, the other one of Player B.

In this interaction Player A will start with 100 MU and Player B with 0 MU. Player A will choose how many MU (up to 100) to send to Player B.

At the end of the interaction, Player A's income will be:  $100 - (\text{How much Player A sent})$

Player B's income will be: (How much Player A sent)

### **Screen 8 – DG Player A**

You start with 100 MU.

If you are player A, how much would you like to send to player B (up to 100)?

[Textbox]

### **Screen 9 – UG Instructions**

You will be matched with one other individual. One of you will be assigned the role of Player A, the other one of Player B.

In this interaction Player A will start with 100 MU and Player B with 0 MU. Player A will choose how many MU (up to 100) to offer to Player B

Player B will choose to accept or reject the offered amount. If Player B accepts the offered amount then they will receive the offered amount and Player A will receive the remainder.

If Player B rejects the offered amount, neither player will receive anything.

At the end of the interaction, Player A's income will be:

$100 - (\text{Offer})$ , if Player B accepts

0 if Player B rejects

Player B's income will be:

$(\text{Offer})$ , if Player B accepts

0 if Player B rejects

### **Screen 10 – UG Player A**

You start with 100 MU.

If you are Player A, how many MU would you like to offer to Player B (up to 100)? [Textbox]

### **Screen 11 – UG Player B**

You start with 0 MU.

Your minimum acceptable offer is the smallest offer which you would accept.

If Player A offers you an amount below this, you will reject, if Player A offers you an amount above this, you will accept.

If you are Player B, what is your Minimum Acceptable Offer?

[Textbox]

## ***12. Instructions: Experiment 2 Stage B***

### **Screen 1 – General Instructions**

You will now be asked to make decisions in 3 different interactions. Unlike the interactions you have experienced so far each of these interactions will last for exactly one round.

For each of these interactions, you will be matched anonymously with new individuals that are currently in the experiment. The choices you make in each interaction will not affect subsequent interactions in any way.

Some interactions will include multiple roles (Player A and Player B), you will be asked to make choices for each role.

After you finish all the interactions, one interaction will be chosen randomly to actually be played. If this interaction had different roles, they will be randomly assigned. Only the outcome of the chosen interaction will affect your (and the other individuals') earned MUs. You will be shown your income from the randomly chosen interaction.

Because you do not know which of your decisions will be chosen it is in your interest to treat each decision as if it counts for real MU.

### **Screen 2 – 3PP DG**

You will be matched with two other individuals. One of you will be assigned the role of Player A, one will be assigned the role of Player B and one will be assigned the role of Observer.

In this interaction Player A will start with 100 MU and Player B with 0 MU. Player A will choose how many MU (up to 100) to send to Player B.

The Observer will be endowed with 100 MU. The Observer will be able to assign Player A an amount of Penalty Units (PU) which will depend on Player A's choice. That is, the Observer will be asked to indicate what amount of PU they wish to assign if Player A chooses to Send 0 MU, what amount they wish to assign if the Player chooses to Send 10 MU and so on.

Each penalty unit will cost the observer 1 MU and will cause the player to whom it is assigned to lose 5 MU from their final earnings.

At the end of the interaction, Player A's income will be:  $100 - (\text{How much Player A sent}) - 5 * (\text{Penalty Units Assigned})$

Player B's income will be:  $(\text{How much Player A sent})$

The Observer's income will be:  $100 - (\text{Penalty Units Assigned})$

### **Screen 3 – DG Player A**

You start with 100 MU.

If you are player A, how much would you like to send to player B (up to 50)?  
[Textbox]

#### **Screen 4 – DG Observer**

Observers start with 100 MU. Player A starts with 100 MU. Player B starts with 0 MU. Each penalty unit that observers end up assigning will cost them 1 MU and cause Player A to lose 5 MU from their final earnings.

If Player A sends 0 MU to Player B, how many penalty units do you wish to assign them?  
[Textbox]

If Player A sends 0 MU to Player B, how “socially inappropriate” do you think this action is? [1-7 scale]

If Player A sends 10 MU to Player B, how many penalty units do you wish to assign them?  
[Textbox]

If Player A sends 10 MU to Player B, how “socially inappropriate” do you think this action is? [1-7 scale]

If Player A sends 20 MU to Player B, how many penalty units do you wish to assign them?  
[Textbox]

If Player A sends 20 MU to Player B, how “socially inappropriate” do you think this action is? [1-7 scale]

If Player A sends 30 MU to Player B, how many penalty units do you wish to assign them?  
[Textbox]

If Player A sends 30 MU to Player B, how “socially inappropriate” do you think this action is? [1-7 scale]

If Player A sends 40 MU to Player B, how many penalty units do you wish to assign them?  
[Textbox]

If Player A sends 40 MU to Player B, how “socially inappropriate” do you think this action is? [1-7 scale]

If Player A sends 50 MU to Player B, how many penalty units do you wish to assign them?  
[Textbox]

If Player A sends 50 MU to Player B, how “socially inappropriate” do you think this action is? [1-7 scale]

#### **Screen 5 – 3PPPD Instructions**

You will be matched with three other individuals. Two of you will be assigned the role of Player and two of you will be assigned the role of Observer.

In this interaction each of the Players will choose an action A or B. If both Players choose A, they both receive 80 MU. If one Player chooses A and the other chooses B the one who chose B will receive 120 MU and the one who chose A will receive 0 MU. If both Players choose B, they each receive 20 MU.



Each Observer will start with 100 MU and will be paired with one Player. The Observer will be able to assign the Player they are paired with Penalty Units or Bonus Units depending on that Player's choices in the game.

Each Penalty Unit will cost the Observer 1 MU and will cause the Player to whom it is assigned to lose 5 MU from their final earnings. Each Bonus Unit will cost the Observer 1 MU and will add 5 MU to the Player's final earnings.

Observers will be able to assign between 0 and 20 Penalty Units and/or between 0 and 20 Bonus Units to the Player they are paired with.

At the end of the interaction, the Players' incomes will be:

If Both Choose A:  $80 - 5 * (\text{Penalty Units Assigned}) + 5 * (\text{Bonus Units Assigned})$

If Both Choose B:  $20 - 5 * (\text{Penalty Units Assigned}) + 5 * (\text{Bonus Units Assigned})$

If the Player chooses A and their opponent chooses B:  $0 - 5 * (\text{Penalty Units Assigned}) + 5 * (\text{Bonus Units Assigned})$

If the Player chooses B and their opponent chooses A:  $120 - 5 * (\text{Penalty Units Assigned}) + 5 * (\text{Bonus Units Assigned})$

The Observer's income will be:  $100 - (\text{Penalty Units Assigned}) - (\text{Bonus Units Assigned})$

### Screen 6 – 3PPPD Player Screen

	A	B
A	80,80	0,120
B	120,0	20,20

If you are a player, which action do you wish to choose? [Button=A] [Button=B]

### Screen 7 – 3PPPD Observer Screen

Observers start with 100 MU and can assign Penalty Units or Bonus Units to the Player they are paired with. Each Penalty Unit causes the Player to lose 5 MU Each Unit from their final earnings. Each Bonus Unit causes the Player to gain 5 MU in their final earnings. Each Unit costs the Observer 1 MU. Observers can assign up to 20 Units.

If the Player you are paired with chooses A, how many Penalty Units do you wish to assign them?

If the Player you are paired with chooses A, how “socially inappropriate” do you think that choice is?

If the Player you are paired with chooses A, how many Bonus Units do you wish to assign them?

If the Player you are paired with chooses A, how “socially praiseworthy” do you think that choice is?

If the Player you are paired with chooses B, how many Penalty Units do you wish to assign them?

If the Player you are paired with chooses B, how “socially inappropriate” do you think that choice is?

If the Player you are paired with chooses B, how many Bonus Units do you wish to assign them?  
If the Player you are paired with chooses B, how “socially praiseworthy” do you think that choice is?

### Screen 8 – PGP Instructions

You are matched with a group of 3 other individuals. Each person in your group is endowed with 100 units. Each individual decides how many of the 100 units they are going to contribute to a common project that benefits all group members (as described below), and how many of them to keep for themselves.

Your group of 4 is also split into pairs of individuals. In addition to choosing how much to contribute to the project, each individual in the group can choose to assign up to 20 Penalty Units to the individual they are paired with, depending on that individual's chosen contribution.

Each Penalty Unit will cost 1 MU and will cause the Player to whom it is assigned to lose 5 MU from their final earnings.

The contributions of all 4 players are added up. The total sum is multiplied by 1.6 and then evenly split among all 4 players. Each player gets the same share from the project.

In addition to your earnings from the project, you also receive the units you chose not to contribute.

Thus, your income is:

$$100 - (\text{your contribution to the project}) + 1.6 \times (\text{sum of all contributions}) / 4 - (\text{Penalty Units You Assign}) - 5 \times (\text{Penalty Units Assigned to You})$$

#### Example

Suppose each person contributes 100 MU to the project. This 400 MU is multiplied by 1.6 and then split between all 4 people.

In this case, each person will receive 160 MU from the project and lose 5 MU for each Penalty Unit assigned to them.

#### Example

Suppose 3 people contribute 100 MU and 1 person contributes 0 MU. The 3 people that contributed will get 120 MU and lose 5 MU for each Penalty Unit assigned to them.

The person that did not contribute will get  
120 (from project) + 100 (that they kept) = 220 MU  
and will lose 5 MU for each Penalty Unit assigned to them.

### Screen 9 – PGP Contribution

You start with 100 MU. You may contribute multiples of (i.e. 0, 25, 50, 75 or 100 MU) to the project.

How much will you contribute?

[Button = 0], [Button=25], [Button=50], [Button=75], [Button=100]

### **Screen 10 – PGP Punishment**

You are paired with one other Player from your group. You may decide how many Penalty Units you wish to assign them based on their actions in the game. Each Penalty Unit will cost 1 MU and will cause the individual to whom it is assigned to lose 5 MU from their final earnings.

If the Player you are paired with contributes 0, how many Penalty Units do you wish to assign them? [Textbox]

If the Player you are paired with contributes 0, how “socially inappropriate” do you think that choice is? [1-7 scale]

If the Player you are paired with contributes 25, how many Penalty Units do you wish to assign them? [Textbox]

If the Player you are paired with contributes 25, how “socially inappropriate” do you think that choice is? [1-7 scale]

If the Player you are paired with contributes 50, how many Penalty Units do you wish to assign them? [Textbox]

If the Player you are paired with contributes 50, how “socially inappropriate” do you think that choice is? [1-7 scale]

If the Player you are paired with contributes 75, how many Penalty Units do you wish to assign them? [Textbox]

If the Player you are paired with contributes 75, how “socially inappropriate” do you think that choice is? [1-7 scale]

If the Player you are paired with contributes 100, how many Penalty Units do you wish to assign them? [Textbox]

If the Player you are paired with contributes 100, how “socially inappropriate” do you think that choice is? [1-7 scale]

### 13. Instructions: Supplemental Experiment

#### RPD Instructions

Thanks for participating!

In this part of the study you will play several *games* with **other people**. **You will play with REAL people also recruited from Mechanical Turk**. Each game will be made of several *rounds*.

#### Payment

- 1) During each game you will earn points.
- 2) For every 10 points you earn, you will be paid an extra 1 cent bonus.
- 3) This bonus is in addition to your payment for completing the HIT.

*Each round will work like this:*

- 1) You and the other person will each make one decision.
- 2) Based on the choices made by **both of you**, you will both receive points.

*Your choices are:*

There are two possible choices available to **both of you**: **A** or **B**.

When a person chooses A, they give the **other person** 50 points.

When a person chooses B, they give **themselves** 20 points.

Option	You get	The other person gets
<b><u>A</u></b> :	0 points	50 points
<b><u>B</u></b> :	20 points	0 points

Your total payoff will be the amount you get from your decision **plus** the amount you get from the other person's decision.

#### Summary

If you **both choose A**, you will both receive 50 points.

If **you choose A** and the **other chooses B**, you will receive 0 points and the other will receive 70 points.

If **you choose B** and the **other chooses A**, you will receive 70 points and the other will receive 0 points.

If you **both choose B**, you will both receive 20 points.

#### How Many Rounds Will There Be?

- 1) At the end of each round, the computer will randomly select either STOP or KEEP GOING.
- 2) If the computer chooses STOP, the current game will end and you will play a new game with a new person.
- 3) If the computer chooses KEEP GOING, you will play another round of the same game with the same person.

The computer will choose KEEP GOING 80% of the time and STOP 20% of the time.

### Repeated PD Comprehension Questions

Which ACTION BY YOU gives YOU a higher payoff in a round? (A or B)

Which ACTION BY YOU gives THE OTHER PLAYER a higher payoff in a round? (A or B)

Which ACTION BY THE OTHER PLAYER gives YOU a higher payoff in a round? (A or B)

Which ACTION BY THE OTHER PLAYER gives THE OTHER PLAYER a higher payoff in a round? (A or B)

### Strategy Choice Condition

#### How You Play:

To play the big games, you will choose a strategy. The computer will automatically makes choices according to your strategy in all of your games.

Remember, in each round you and the other person make one choice:

Choosing A gives the other 50 points.

Choosing B gives you 20 points.

<u>Option</u>	<u>You get</u>	<u>The other person gets</u>
<u>A:</u>	0 points	50 points
<u>B:</u>	20 points	0 points

You can choose between two strategies:

**Strategy 1** - This strategy will start every big game by choosing A. It will see what the other person does on that round.

If the **other person chooses A in a round**, this strategy will choose **A in the next round**.

If the **other person chooses B in a round**, this strategy will choose **B in the next round**.

**Strategy 2** - This strategy will always choose **B in every round**.

Please use the buttons below to indicate which strategy you would like to pick.

#### (Active Play) Repeated PD New Person

You are now matched with a NEW person for a NEW game!

#### (Active Play) Repeated PD Round Choice Screen

You and the other person will now each choose A or B.

Remember, for each of you:

**Choosing A gives the other 50 points.**  
**Choosing B gives you 20 points.**

<u>Option</u>	<u>You get</u>	<u>The other person gets</u>
<u>A</u> :	0 points	50 points
<u>B</u> :	20 points	0 points

**Your choice: (A or B)**

**(Active Play) Repeated PD Results Screen**

You chose:  $X$

The other person chose:  $Y$

This means you earn  $L$  points and the other person earns  $K$  points.

[The computer has chosen KEEP GOING, there will be another round.] OR [The computer has chosen STOP, this was the last round of this game]

**DG Instructions**

You will now play a new game with a new person also recruited from Mechanical Turk.

In this new game, you are given an extra 80 cents. You then choose how to split this money between yourself and the other person.

The other person did not play any other games during this study, and cannot affect your earnings. The only bonus they receive is the money you give them.

You have now been given an extra 80 cents. You may choose how to split this money between yourself and the other person.

Use the buttons below to indicate how much of this 80 cents you would like to transfer to the other person. You keep all the remaining money (for example: if you transfer 10 cents, you will keep 70 cents).

Please choose how many cents you will transfer to the other person:  
 (Radio buttons 0-40 in increments of 10)