



*The World's Largest Open Access Agricultural & Applied Economics Digital Library*

**This document is discoverable and free to researchers across the globe due to the work of AgEcon Search.**

**Help ensure our sustainability.**

Give to AgEcon Search

AgEcon Search

<http://ageconsearch.umn.edu>

[aesearch@umn.edu](mailto:aesearch@umn.edu)

*Papers downloaded from **AgEcon Search** may be used for non-commercial purposes and personal study only. No other use, including posting to another Internet site, is permitted without permission from the copyright owner (not AgEcon Search), or as allowed under the provisions of Fair Use, U.S. Copyright Act, Title 17 U.S.C.*

*No endorsement of AgEcon Search or its fundraising activities by the author(s) of the following work or their employer(s) is intended or implied.*

**Cooperation in Social Dilemmas with Correlated Noisy Payoffs: Theory and  
Experimental Evidence**

**Alecia Evans<sup>1</sup>, Department of Agricultural Economics, Purdue University,  
evans264@purdue.edu**

**Juan Sesmero, Department of Agricultural Economics, Purdue University,  
jsesmero@purdue.edu**

***Selected Paper prepared for presentation at the 2021 Agricultural & Applied Economics Association  
Annual Meeting, Austin, TX; August 1-August 3***

*Copyright 2021 by Evans and Sesmero. All rights reserved. Readers may make verbatim copies of this document for non-commercial purposes by any means, provided that this copyright notice appears on all such copies.*

---

<sup>1</sup> Corresponding author

# COOPERATION IN SOCIAL DILEMMAS WITH CORRELATED NOISY PAYOFFS: THEORY AND EXPERIMENTAL EVIDENCE

Alecia Evans and Juan Sesmero

## **Abstract**

In infinitely repeated social dilemmas, forces that cloud knowledge about past behavior may induce subjects to incorrectly infer their opponents' past actions, possibly inhibiting cooperation. Random shocks affecting players' payoffs constitute one such force. We develop a framework to study this environment and predict that correlation across shocks can restore cooperation by enhancing knowledge about past behavior. We then test this prediction in a laboratory experiment. On average, we fail to confirm our prediction. Nevertheless, we find that correlation across shocks fosters (inhibits) cooperation among subjects that choose to cooperate (defect) during the initial stages of the game. We complement our experiments with simulations based on a genetic algorithm and find that correlation makes conditional cooperation strategies more successful, prompting these strategies to survive the evolutionary process. As a result, in an evolutionary framework, correlation unambiguously enhances cooperation.

JEL Codes: C73; C92; D81; D82

# 1 Introduction

Many important economic activities are carried out in groups where agents interact repeatedly over time. These groups are often formed to overcome market failures that inhibit socially desirable trading. For instance, agents organize in groups to facilitate informal risk-sharing (e.g. Fitzsimons, Malde, and Vera-Hernández (2018)), provide access to information, insurance and credit (e.g. Bloch, Genicot, and Ray (2008)), improve productivity and profitability of farmers (Agarwal 2018), and prevent resource exhaustion (e.g. Ostrom et al. (1999)). The success of these groups, however, crucially depends on the ability of members of the group to cooperate with each other.

But many of these settings have the structure of a social dilemma – a situation in which rational agents may fail to cooperate even when cooperation is mutually beneficial. This problem is exacerbated by uncertainty about past behavior – for example, arising from agents’ inability to perfectly monitor each other – which may induce subjects to incorrectly infer others’ past actions. Imperfect monitoring can be induced by random shocks that alter agents’ payoffs (Bendor, Kramer, and Stout 1991; Bendor 1993). Such environments are empirically pervasive. As such, in this paper we focus on imperfect (private) monitoring and investigate how the structure of correlation across random shocks that alter players’ payoffs affects the strength of monitoring and, ultimately, cooperation in infinitely repeated social dilemmas.

The outcome of infinitely repeated social dilemmas greatly depends on the strategic environment (see Dal Bó and Fréchette (2018) for a recent review). An important distinction between the structure of correlation across shocks and other features of the trading environment that can spur or hinder cooperation, is that correlation also affects risk-sharing; and risk-sharing is at the core of many economic activities that are carried out in groups. In developing countries, where formal insurance and credit markets are underdeveloped and economic well-being is very sensitive to random shocks (including income and health shocks), informal risk-sharing arrangements serve as a safety net. A prominent example of this environment is the extended family. In this context, risk-sharing takes the form of reciprocal credit systems where some siblings assist others who later reciprocate (Baland et al. 2016), or a buffer to smooth consumption in the event of crop losses (Fitzsimons, Malde, and Vera-Hernández 2018). In these settings, negative correlation across shocks facilitates risk-sharing while positive correlation inhibits risk-sharing (e.g. Fafchamps (2011)).

But in addition to the correlation structure facilitating risk-sharing, cooperation is necessary for the success of risk-sharing groups. Yet, cooperation often fails in these environments. Members of the family sometimes avoid sharing wealth by taking out loans to feign liquidity constraints (Baland, Guirking, and Mali 2011), or try to obscure their true endowments from others in the family (Jakiela and Ozier 2016). When agents interact repeatedly, cooperation is more likely, but hardly a foregone conclusion. While the effect of the correlation structure on risk-sharing is straightforward, its effect on cooperation remains unclear.

The primary question we raise in this study is whether the structure of the correlation across noisy payoffs affects cooperation among subjects. But we are also interested in understanding the mechanisms underlying this effect. One possibility is that correlation strengthens monitoring, thereby allowing subjects to lower inferential error. This lower inferential error could also prompt a change in

the strategies used by players.

We investigate these issues in three steps. First, we develop a theoretical framework and generate testable predictions regarding the effect of correlation on inferential error and, ultimately, cooperation. We then test these predictions in a laboratory experiment and examine other mechanisms arising in the experimental setting. Finally, we complement the experiment with simulations based on a genetic algorithm. This allows us to examine which strategies observed in the experiment are likely to survive from an evolutionary point of view.

We build on Bendor, Kramer, and Stout (1991) and Bendor (1993) and develop a framework to formally model behavior in an infinitely repeated prisoner’s dilemma with noisy payoffs. We extend this framework by 1) allowing players to choose the benchmark against which a private signal is defined, and 2) allowing shocks affecting subjects’ payoffs to be (positively and negatively) correlated. Based on this framework we predict that correlation will strengthen monitoring (that is, will lower inferential errors by subjects) thereby enhancing cooperation.

We test the predictions generated from our theoretical framework in a laboratory experiment. A laboratory experiment is appropriate because the private information needed to understand how monitoring impacts cooperation is not usually available from observational data. Additionally, and in contrast to field experiments, a laboratory experiment allows us to have full control of the strategic environment. We can exogenously manipulate the correlation structure and prevent communication outside of the strategic environment which allows us to establish clear causality and identify subjects’ inferences and strategies within a large set of possible options.

In our experiment, subjects play an infinitely repeated prisoner’s dilemma (PD) with a continuation probability of  $\delta = 0.9$ . In the stage game of the PD, an agent’s payoff is affected by a random shock that is uniformly distributed with mean zero. Each player receives a private noisy signal about the realized payoff of the other player in relationship to a benchmark value set by the subject. The subject can combine this information, with information on her own payoff, to infer whether the other agent has deviated or defected. We then expand this by including correlated shocks. We implement four treatments that vary the correlation level from high ( $\rho = 0.9, -0.9$ ) to moderate ( $\rho = 0.4, -0.4$ ), to compare against a baseline case of no correlation ( $\rho = 0$ ).

Our main result is that correlation, either positive or negative, does not improve cooperation relative to the baseline of  $\rho = 0$ , on average. A closer look at mechanisms clarifies this seemingly puzzling result. Stronger correlation does tend to lower inferential error. But this does not persuade “not nice” subjects (not nice subjects are those that defect in the very first round) to engage in cooperation, unconditional or otherwise. Therefore, in games where a “not nice” player is involved, correlation helps unveil defection which precipitates the unraveling of cooperation. This is confirmed by experimental results which show that, when both players are “not nice”, correlation is associated with lower inferential error and lower cooperation. Conversely, and by the same mechanism, when both players are “nice”, correlation is associated with lower inferential error and higher cooperation.

Across observations, most of the interactions involve “not nice” players. This fact underpins the muddled relationship we find between correlation and cooperation, on average. We complement our laboratory experiment with a computational experiment based on an evolutionary algorithm to

provide further intuition on how cooperation can be maintained overtime under such environment. In this process, certain types of players, those implementing the most successful strategies, are more likely to survive (Axelrod 1980). Indeed, we find that higher degrees of correlation did result in high levels of cooperation.

The rest of the paper is organized as follows. In Section 2, we discuss the nature of our contribution in the context of the broader literature. In Section 3, we present the theoretical background. In Section 4, we give the details of the experimental design. In Sections 5 and 6, we outline the questions our analysis will answer and the main results from our experiment. In Section 7, we present a computational experiment that test behavioral aspects of our experimental design. In Section 8, we conclude with a discussion of our main results.

## 2 Related Literature

We contribute to the literature on cooperation in infinitely repeated PDs when knowledge about past behavior is limited, leading to inferential uncertainties about other players' actions (players are forced to guess their opponents' past behavior). A part of this literature introduces inferential uncertainty by considering noise in the form of implementation error, experimentally and theoretically (Fudenberg and Maskin 1990; Miller 1996; Fudenberg, Rand, and Dreber 2012; Imhof, Fudenberg, and Nowak 2007; Ioannou 2014a, 2014b; Zhang 2018). With implementation errors, there is a probability that the action the players implement is different from the one they intended. This obscures knowledge about past behavior in the sense that players know the action their opponent took but are unsure about their intentions. Furthermore, players are aware of this probability. We can consider this as a signal each player receives about the probability that their opponent actually intended the observed action. This signal delivers information, albeit incomplete. Consequently, players may incorrectly infer the intent of others. Papers in this strand of literature find that incomplete information regarding intent can, though not always, reduce cooperation.

A key feature of the literature on implementation error is that the environment is characterized by imperfect information (where past actions are observable, but the intention is unclear) rather than imperfect monitoring (where past actions are unobservable). This is a subtle, yet important distinction. Both frameworks are appropriate for distinct empirical settings; and they are not observationally equivalent, that is, one does not tend to mimic the other. As pointed out by Ioannou (2014b), imperfect monitoring (which causes individuals to draw incorrect inferences about others past actions) is more detrimental than imperfect information (where individuals are prone to errors in implementing their own actions). This is because implementation errors can introduce cooperative actions even in the presence of unconditional uncooperative strategies (for example, Always Defect), thereby facilitating cooperation. While this would also imply possible deviations from unconditional cooperation strategies (thereby hindering cooperation), these kinds of strategies are not as ubiquitous as their unconditional defection counterparts. In this study, we employ an imperfect monitoring framework because it better captures key features of the empirical settings that motivate our analysis.

We examine the literature on imperfect monitoring in two broad strands. One strand of the literature on imperfect monitoring studies deterministic PDs where the players do not directly observe

their opponent’s past action, but receive a private signal with a set accuracy about such action (Aoyagi, Bhaskar, and Fréchet 2019; Kayaba, Matsushima, and Toyama 2020). The signal is either good or bad, and a good signal is more likely to occur when their opponent is cooperative. The monitoring accuracy is the probability of receiving the correct signal, and a lower accuracy translates into a higher probability of inferential error (that is, a higher chance that a player will incorrectly guess their opponent’s past action). Aoyagi, Bhaskar, and Fréchet (2019) vary the monitoring environment and find that subjects can sustain cooperation under imperfect private monitoring (at rates comparable to perfect monitoring but lower than imperfect public monitoring). Kayaba, Matsushima, and Toyama (2020) vary the accuracy of the signal and find that cooperation increases as monitoring strengthens, that is, as the signal becomes more accurate.

The empirical settings that motivate our study, such as the extended family, collective agrarian societies, and micro-finance groups, involve random shocks that affect subjects’ payoffs (for example, weather events, unexpected health issues). Therefore, while the deterministic PD framework employed by Aoyagi, Bhaskar, and Fréchet (2019) and Kayaba, Matsushima, and Toyama (2020) captures the key issue of imperfect monitoring, it does not fit situations where random shocks affecting payoffs constitute the source of imperfect monitoring. The other strand of literature on imperfect monitoring introduces uncertainty regarding past behavior through noise in the form of random payoffs (Bendor, Kramer, and Stout 1991; Bendor 1993). We build on the framework developed by Bendor, Kramer, and Stout (1991) and Bendor (1993), but our analysis differs from those papers in important ways.

First, our primary objective is to understand how correlation across shocks affecting payoffs alters inferential error (monitoring strength) and, ultimately cooperation. We study this because in the empirical settings where payoffs are affected by random shocks, these shocks are often correlated. In many cases, shocks are positively correlated. For instance, in many microfinance institutions, in order to overcome moral hazard and adverse selection, groups are composed of agents living in the same geographic space and probably conducting similar economic activities such as farming. In other cases, shocks are negatively correlated. For instance, in many extended family settings (informal insurance), players engage in fundamentally different activities such as farming and urban employment; activities that are often negatively correlated, or uncorrelated. To better understand cooperation in these settings, we extend the framework in Bendor, Kramer, and Stout (1991) and Bendor (1993) to accommodate correlation across shocks. We then systematically vary the correlation structure and compare inferential errors and cooperation across structures.

Like in Bendor, Kramer, and Stout (1991) and Bendor (1993), players in our framework receive a signal about their opponent’s payoff. The signal is defined in relationship to a benchmark value, that is, the signal indicates whether the opponent’s payoff is above or below that benchmark. A salient feature of our framework is that we allow subjects in the lab (and automata in the genetic algorithm) to choose the benchmark. The informational value of the signal (the extent to which the signal helps players infer their opponent’s action) depends upon where the benchmark is set. Moreover, the correlation structure affects the informational content of the signal, but the degree to which this information is exploited depends on, once more, where the benchmark is set. Therefore, in our framework, the accuracy of the signal is endogenous – it depends upon the subjects’ ability to set the benchmark at a level that minimizes inferential error. We now turn to a more formal characterization of our framework.

### 3 Theoretical Background

We study a situation in which two players play an infinitely repeated prisoner’s dilemma game. The deterministic payoff is that of a standard prisoner’s dilemma (see Table 1) where,  $T > R > P > S$  and  $2R > T + S$ . However, for the stage game, the payoff to each player is affected by a uniformly distributed shock. With this shock, the realized payoff (stage game payoff plus random shock) becomes  $\hat{T}$ ,  $\hat{R}$ ,  $\hat{P}$  and  $\hat{S}$ . This payoff is  $\hat{X} = X + V$ , where  $X = \{T, R, P, S\}$ , and  $V$  is the random shock with mean zero, and it is uniformly distributed between a lower bound  $V_{LB}$  and an upper bound  $V_{UB}$ . That is,  $V_{LB} \leq V \leq V_{UB}$ . Two shocks are generated, one for each agent, which we denote by  $V_1$  and  $V_2$ . The shocks are independent across rounds but can be positively or negatively correlated between them.

Table 1: Deterministic payoff for the prisoner’s dilemma

|          | <b>C</b> | <b>D</b> |
|----------|----------|----------|
| <b>C</b> | R, R     | S, T     |
| <b>D</b> | T, S     | P, P     |

Similar to Bendor (1993), for a range of payoffs, the random shock introduces uncertainty and limits a player’s ability to infer their opponent’s actions. In this range of payoffs, which we call the region of uncertainty, if a player plays cooperate, it is possible to incorrectly infer that the other player defected when they had in fact cooperated (Type 1 error). Also, they could incorrectly infer that the other player had cooperated when they had in fact defected (Type 2 error). This region of uncertainty exists as long as  $S + V_{UB} > R + V_{LB}$ .<sup>1</sup>

Figure 1 gives the distribution of the realized payoff of player 1 when she cooperates. For any realized payoff to the left of the region of uncertainty, player 1 knows without a doubt that they received the sucker payoff,  $\hat{S}$  (the payoff is too low to be anything else). And for any region to the right, player 1 knows without a doubt that they received the reward payoff,  $\hat{R}$  (conditional on the subject having cooperated, the payoff is too high to be anything else). Within the region of uncertainty, player 1 is unsure and there is a probability  $p > 0$  that they will make an incorrect inference about the other player’s action. We call this an inferential error. The reason for player 1’s uncertainty within this range of payoffs is that two scenarios are probabilistically possible. It is possible that player 2 defected, but that player 1 received a large and positive shock, making the payoff that of a “lucky sucker”. It is also possible that player 2 cooperated, but that player 1 received a large and negative shock, making the payoff that of an “unlucky reward”.

We assume that shocks are uniformly distributed. This departs from the treatment in Bendor (1993) which assumed that shocks are normally distributed (see Appendix A). To better fit the empirical settings that are of primary interest to us, we introduce correlation across these shocks. This also departs from Bendor (1993) who assumed that shocks are independent across players. We do, however, maintain that shocks are uncorrelated over time. As previously established, when people are engaged in social dilemma type groups, correlation among shocks is the norm, rather than the

1. If the player plays defect, the region of uncertainty exists as long as  $P + V_{UB} > T + V_{LB}$



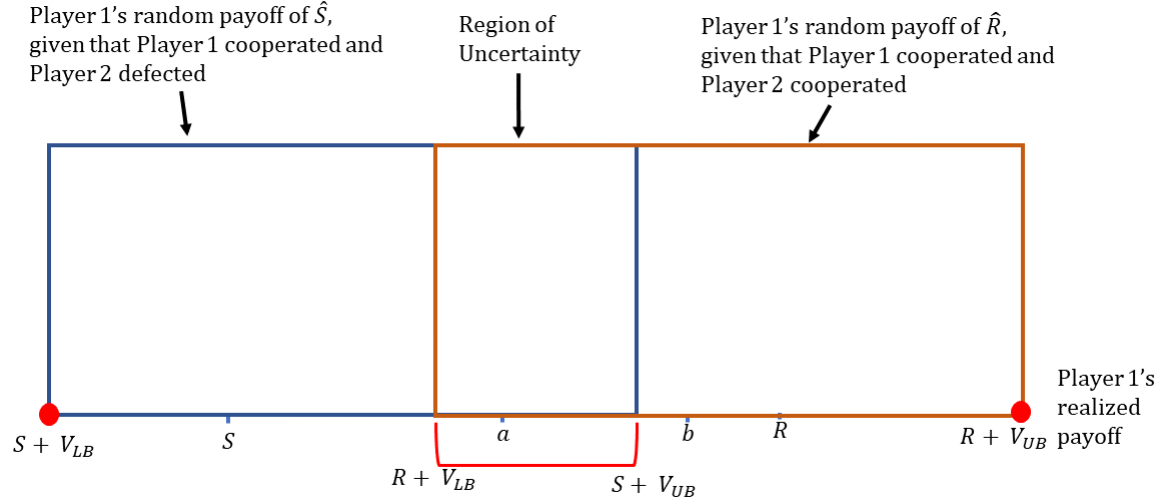


Figure 1: Player 1's realized payoff when cooperation is played. The height of the distribution is the probability of the realized payoff  $\hat{X}$

exception. As such, we conjecture that knowledge of the correlation structure may help reduce a players' inferential error. In particular, we theorize that knowledge of the correlation structure can reduce inferential error by allowing for some degree of monitoring.

To see how this happens, imagine that player 1 cooperates and receives a payoff slightly above  $R + V_{LB}$  (see Figure 1). That is, she receives a payoff that is low within the region of uncertainty. Because player 1 is in the region of uncertainty, without any other information, it is hard for player 1 to know whether player 2 cooperated and she got unlucky (received a bad shock), or whether player 2 defected and she got lucky (received a good shock). Now, let us assume that player 1 knows that player 2 received a payoff above some benchmark value (Bendor's critical cutoff value from Appendix A). For the sake of argument, let this benchmark value be greater than  $R$  in Figure 1. This indicates that player 2 did well in that they received a high payoff. This could have resulted from two likely situations: (1) player 2 defected or (2) player 2 cooperated and she simply got lucky while player 1 did not. But, if player 1 knows that shocks are positively correlated across players, then a scenario in which player 2 defected is more likely than one in which player 2 cooperated and received a bad shock, while player 1 received a good shock.

We formalize this as follows. Consider two states of nature  $\theta \in \{C, D\}$ , indicating that the other player has cooperated ( $C$ ) or defected ( $D$ ). Each player starts with an uninformed prior that these events are equally likely. Then, each player receives a signal,  $s \in \{0, 1\}$  that allows them to update their belief using simple Bayesian techniques. The signal tells a player that the other player's payoff is above or equal to ( $s = 0$ ) or below it ( $s = 1$ ) a benchmark chosen by her. The signal is the imperfect information the player has regarding the payoff of the other.

We now formalize a decision rule using the noisy signal. Assume that player 1 gets a signal that tells her where player 2's realized payoff lies in relation to this benchmark value. She then uses this signal to update her belief on  $C$  and  $D$ . The simple Bayesian process is outlined in Appendix B.

The probability that the opponent will be above the benchmark, conditional on the player having cooperated,  $P(s = 0/\theta = C)$ , is denoted by  $\pi_C$ . Similarly, the probability that the opponent will be below the benchmark, conditional on the player having defected,  $P(s = 1/\theta = D)$ , is denoted by  $\pi_D$ . The signal player 1 receives about player 2's realized payoff gives additional support about  $\theta = C$  if  $\pi_C > 1 - \pi_D$ . Due to symmetry, a similar argument holds for player 2.

Notice that  $s = 0$  is equivalent to  $P > B$ , where  $P$  is the other player's realized payoff and  $B$  is the benchmark chosen. Similarly,  $s = 1$  is equivalent to  $P < B$ . Therefore, the expression  $\pi_C > 1 - \pi_D$  is actually a function of the benchmark chosen by the player. If the player sets a very high benchmark (a value close to the upper bound of the possible realized payoffs of the other player), then  $\pi_C < 1 - \pi_D$ . In this case, the player will always infer defection regardless of the realized payoff of the other player. In turn, this implies high levels of Type I error and low levels of Type 2 error. Conversely, if the player sets a very low benchmark (a value close to the lower bound of the possible realized payoffs of the other player), then  $\pi_C > 1 - \pi_D$  regardless of the other player's realized payoff. Therefore, she will always infer cooperation by her opponent, which in turn implies high levels of Type 2 error and low levels of Type I error. In both of these cases, inference will be incorrect close to half of the time. Therefore, if the benchmark is set close to the upper or lower bound of the other player's possible realized payoffs, the signal regarding the other player's realized payoff contains very little information and may not lead to better inference.

On the other hand, if the player sets a more intermediate level for the benchmark, then Type I and 2 errors will be balanced (similarly frequent), and overall errors will be minimized. But the level will depend on how much of an overlap there is with these two distributions. These distributions and consequently their overlaps are shaped by both the correlation between shocks and the actions of player 1. This is illustrated numerically in Appendix B and C. In Figures 6 and 7 we simulated the realized payoff for player 2, for different combinations of player 1's realized payoff (within the region of uncertainty) and correlation between shocks, conditional on player 1 cooperating.

As shown by the numerical simulations in Figures 6 and 7 (Appendix C), a higher correlation (either positive or negative) between shocks shifts the distributions apart, thereby reducing inferential errors, all else constant. This is a mechanical effect and is independent of where the player sets the benchmark based on which the signal is defined. Wherever the benchmark is set, if player 1 infers defection (cooperation) when player 2's payoff is above (below) the benchmark, higher correlation between shocks will make it more likely to that this is in fact true.

But the degree to which correlation translates into a reduction in inferential errors also depends upon where player 1 sets the benchmark. If player 1 sets the benchmark at the payoff where both distributions intersect and correlation is high, the signal will convey information highly indicative of player 2's actions which results in low inferential error and, moreover, a situation where Type I and 2 errors are equally likely. In other words, by choosing the correct benchmark, player 1 can, with a given signal, refine her Bayesian updating of the prior inference regarding player 2's actions. Therefore, there is also a behavioral channel through which higher correlation reduces inferential error. Higher correlation reduces inferential error the most when the agent has the ability to choose the right benchmark based on which the signal is defined.

We illustrate the effect of correlation on inferential error in Figures 8 and 9 in Appendix C,

where we present the probability distributions of player 2’s realized payoffs. We assume that, after the signal, Bayesian updating proceeds based on the benchmark value that equates Type 1 and Type 2 errors (from Player 2’s distribution) when  $\rho = 0$ . In Table 14 in Appendix B we present Type 1 and Type 2 errors along with  $\pi_C$  and  $\pi_D$  using these benchmark values. Using Bayesian updating we put forward a simple decision rule for player 1. If  $\rho \geq 0$  and player 2’s realized payoff is above the benchmark value, assume that player 2’s mostly likely action was defection. Likewise, if player 1 is signaled that player 2’s realized payoff is below the benchmark value, assume that their most likely action was cooperation. The opposite holds for  $\rho < 0$ . Simulations reported in Table 14 show that higher correlation, under this choice of benchmark, translates into a significant reduction of both type I and 2 errors, but the reduction is larger when correlation is positive.

We hypothesize that players in an infinitely repeated prisoner’s dilemma will use the signal more effectively when correlation is higher, that this will strengthen monitoring between players (by reducing inferential error), and that this will in turn enhance cooperation. To test these hypotheses we implemented an experiment, which we now proceed to discuss.

## 4 Experimental Design

The experiment is designed to test if a stronger correlation between shocks affecting players’ payoffs reduces inferential error and by extension fosters cooperation in an infinitely repeated prisoner’s dilemma. To induce the infinitely repeated game, subjects were informed that after each round, there was 0.9 probability that a supergame will continue for another round. We pre-drew the random game length of each supergame to ensure that in each session, each supergame lasted for same number of rounds. For each treatment, subjects played 84 rounds over 10 supergames (See Table 2 for the treatment summary).<sup>2</sup>

Table 2: Treatment, sessions, and subjects in the experiment

| Treatment     | No. of Sessions | Total Subjects |
|---------------|-----------------|----------------|
| $\rho = 0.9$  | 3               | 34             |
| $\rho = 0.4$  | 3               | 36             |
| $\rho = 0$    | 3               | 36             |
| $\rho = -0.4$ | 3               | 36             |
| $\rho = -0.9$ | 3               | 36             |

Table 3 shows the stage game of the prisoner’s dilemma, denoted in points. In each round, each subject’s payoff was affected by random, uniformly distributed shock in the range  $[-24, 24]$ . Subjects were only told their realized payoff (payoff inclusive of random shock faced) and the correlation level between their random shock and the random shock of the other player. All these details were included in the instructions that the subjects read on their computer monitors at the beginning of each session. The subjects also had access to the same instructions in written form throughout the entire session. An example of these instructions is included in the Appendix D. We implemented five treatments: a

2. Across all 5 treatments, subjects made a total of 14,952 decisions.

very high positive and negative correlation ( $\rho = 0.9$  and  $\rho = -0.9$ ), moderate positive and negative correlation ( $\rho = 0.4$  and  $\rho = -0.4$ ) and the baseline case of no correlation ( $\rho = 0$ ).

Table 3: Payoff of the stage game

|   | C      | D      |
|---|--------|--------|
| C | 48, 48 | 13, 60 |
| D | 60, 13 | 25, 25 |

Before each supergame, each subject selects two benchmark values, one to be used if they select cooperate and the other if they select defection. The benchmark values were restricted to the range of possible realized payoffs of their opponent. That is, subjects were restricted to the range [24, 84] if they choose cooperation and a range of [-11, 49] if they choose defection. To assist subjects in understanding how the benchmark values work, we included an interactive feature in the instructions. This is a simulation in which the subject and the computer simultaneously make a choice, then the subject receives a feedback on their realized payoff. Also, there is a slider that allows subjects to play around with setting different benchmark values to see what feedback they will receive (above, below, or equal to) about the other player’s realized payoff.

At the beginning of each round, subjects choose between cooperation and defection (in the experiment, we used neutral language of “A” and “B”, instead of “Cooperate” and “Defect”). Immediately after this choice they receive feedback on the resulting realized payoff and also a private signal about the realized payoff of the other subject. This signal tells if the other subject’s realized payoff is above, equal to, or below the benchmark value each subject selected at the beginning of each supergame.

In our design, we opted to allow subjects to select their own benchmark to mimic noisy signals in real-world group interactions. In real-world groups, for example the extended family or other similar settings, individuals use spending habits of others to determine their well-being (Baland, Guirkingier, and Mali 2011; Jakiela and Ozier 2016). Individuals may vary on the thresholds above which they infer defection. For example, an individual may consider that the other person’s paying rent or debts is sufficient proof that they are shirking on risk-sharing agreements. In contrast, others may set a higher bar and consider traveling or similar bigger spending event as sufficient proof of defection.

After each round, we also used a Binarized Scoring Rule (BSR) to elicit incentivized beliefs from each player about the actions of the other player. This was necessary to help us estimate inferential errors and also to estimate the strategies that are being played by subjects across supergames. In each round, after subjects make a decision and their realized payoff is revealed, they were asked how likely they believed that the other subject cooperated. The BSR is incentive compatible in that, as long as subjects prefer getting a reward as opposed to no reward, to maximize the probability of getting the reward, the best action is for them to truthfully report their beliefs about the other subject’s action (Hossain and Okui 2013). In our context, the reward was an additional 2 points. That is, over 84 rounds, a subject could earn a maximum of 164 additional points for truthfully reporting their beliefs about the likely action of the other player. The BSR is also independent of risk attitudes and whether the subject is an expected utility maximizer or not. In Appendix E, we describe the belief elicitation process. In the instructions, we did not give subjects the full details of the belief elicitation process.

They were informed that the details were available after the session.<sup>3</sup> This design feature is consistent with the results of Danz, Vesterlund, and Wilson (2020), who, in an experiment using BSR, found that transparent information on incentives gave rise to error rates in excess of 40%.

On the decision screen for each round, each subject saw a summary of the decision and outcome from the previous round. Also, there was a brief summary of the correlation structure for the treatment. There was also a reminder of the two benchmark values that they selected. At the end of a supergame, before they are randomly rematched, subjects received a detailed account of the actions they made and their realized payoffs, as well as the actions and payoff of the other subject. See Appendix F for a screenshot of this.

A total of 178 subjects participated in 15 sessions at Purdue University’s Vernon Smith Experimental Economics Laboratory (VSEEL). Each treatment had either 10 or 12 subjects and lasted for 90 minutes.<sup>4</sup> Subjects accumulated points during each session, and these were converted at an exchange rate of \$1 = 300 points. On average, subjects earned \$21.58 including a show-up fee of \$5. We used a between subject design, where each subject participated in only one session. For the session, subjects first read the on-screen instructions, then they completed quiz. Following this, they played five unpaid practice rounds against the computer. For additional practice in setting the benchmark, in the practice rounds subjects were allowed to select the benchmark after each round. But for the paid repeated games, the benchmark values were set at the beginning of each supergame. They were informed of all of this in the instructions at the beginning of each session. For the paid repeated prisoner’s dilemma, after each supergame, subjects were randomly rematched with another subject in the room.<sup>5</sup> The experiment was programmed in oTree (Chen, Schonger, and Wickens 2016).

## 5 Questions and Predictions

Based on the theoretical analysis in Section 3, numerical simulations (reported in Appendix B and discussed in Section 3), and insights from previous literature, we predict that a stronger correlation between shocks that affect players’ payoffs, fosters cooperation. We further hypothesize that correlation fosters cooperation by reducing inferential error, that is, by strengthening monitoring. With lower inferential error, we theorize that subjects will depend on more lenient strategies with stronger correlation. We start by stating our primary question:

**Question 1:** *Is cooperation higher with stronger correlation?*

Based on insights from previous studies (Kayaba, Matsushima, and Toyama 2020) we make the following prediction:

**Prediction 1:** Cooperation will increase with stronger positive correlation across shocks. The effect of stronger negative correlation is non-monotonic and, thereby, ambiguous. In other words, stronger positive correlation will foster cooperation, but stronger negative correlation may foster or hinder

---

3. No subject asked for this information after the session.

4. COVID-19 regulations stipulated a maximum of 13 subjects per session. For some sessions, we had no-shows that resulted in only 10 subjects per session.

5. This was not a perfect random rematching. There was a possibility that a subject could be re-matched with someone they played with in a previous supergame.

cooperation.

We now turn our attention to the mechanism underlying Prediction 1. First, we focus on the link between correlation and monitoring strength or, in other words, inferential error:

**Question 2:** *Is inference error lower with stronger correlation?*

We will focus on the degree of inferential error across the various correlation structures. This includes total error rate, as well as Type 1 and Type 2 errors. From our numerical simulations based on Section 3 (see Appendix B), we predict that the correlation structure will affect inferential errors through a mechanical channel and a behavioral channel. We first focus on the mechanical channel, whereby a stronger correlation reduces inferential error even without the subjects using the signal to update their beliefs about the other player’s action. Based on our theoretical analysis, we make the following predictions:

**Prediction 2:** Stronger positive correlation reduces inferential errors, thereby improving the ability of subjects to monitor each other. It does so by lowering Type 1 and 2 errors.

**Prediction 3:** The effect of negative correlation on inferential error is non-monotonic. Weak negative correlation raises inferential error (both Type 1 and 2), but strong negative correlation reduces it (both Type 1 and 2). Therefore, as negative correlation becomes stronger, it first impairs and then improves the ability of subjects to monitor each other.

Predictions 2 and 3 constitute one channel through which correlation affects cooperation in the way described by Prediction 1. This follows from previous studies, which suggest that stronger monitoring (lower inferential error) has a positive effect on cooperation (Kayaba, Matsushima, and Toyama 2020). These findings suggest that positive correlation is likely to foster cooperation because it lowers inferential error. Our simulations suggest that, in turn, negative correlation has an ambiguous effect on cooperation because it has an ambiguous effect on inferential error.

We now discuss the behavioral channel, whereby subjects can use a combination of the correlation structure, and a noisy private signal about their opponent’s payoff to update their beliefs on the likely action of their opponent. The noisy signal is determined relative to a benchmark set by the player. The benchmark is some payoff on the domain of the opponent’s payoffs. As portrayed in Figures 6 and 7, conditional on a player’s action, payoff received, and a correlation structure, there are two payoff distributions for her opponent: one when the opponent cooperates, and one when she defects. These two distributions can overlap; the overlap becomes smaller as the correlation between shocks strengthens and eventually disappears when correlation is sufficiently strong. If the benchmark is set at the center of the overlapping region, Type I and II errors are roughly equalized, and the total inferential error is minimized. Based on the patterns of the overlapping distributions (e.g., Figures 8 and 9), when shocks are positively (negatively) correlated, the opponent is likely to have cooperated (defected) if her payoff is above the threshold. Therefore, we make the following predictions regarding the players’ inferential process:

**Prediction 4:** Subjects will set a benchmark near the center of the overlapping region of the opponent’s payoff distributions.

Then, if the player’s payoff falls within her region of uncertainty, she will use the following decision

rule: when  $\rho \geq 0$ , a signal that the other player is above the benchmark value supports the inference that the other player defected in the previous round. And, when  $\rho < 0$  such a signal supports the inference that the other player cooperated.

In summary, we predict that positive (negative) correlation will reduce (have an ambiguous impact on) inferential errors regardless of the player’s inferential process (Predictions 2 and 3). But we also predict that players will adjust their inferential process to minimize total inferential error and equalize Type I and Type II errors (Prediction 4). If Prediction 4 is correct, the players can more effectively use the information contained in the signal, conditional on the correlation. In other words, correlation is more likely to reduce inferential error if players follow the inferential process described in Prediction 4.

We now look at the effect of correlation on the type of strategies used by subjects, which will influence the prevalence of cooperation. We raise the following question:

**Question 3:** *Do subjects play more lenient strategies when random shocks are correlated?*

Drawing on conjectures in Aoyagi, Bhaskar, and Fréchette (2019), we make the following prediction about dominant strategies under alternative correlation structures:

**Prediction 5:** If a player thinks, with more confidence, that his opponent cooperated, they will be more lenient. This is to reduce the likelihood of retaliation. Combining Predictions 1-4, we expect that strong positive correlation will lead subjects to employ more lenient strategies.

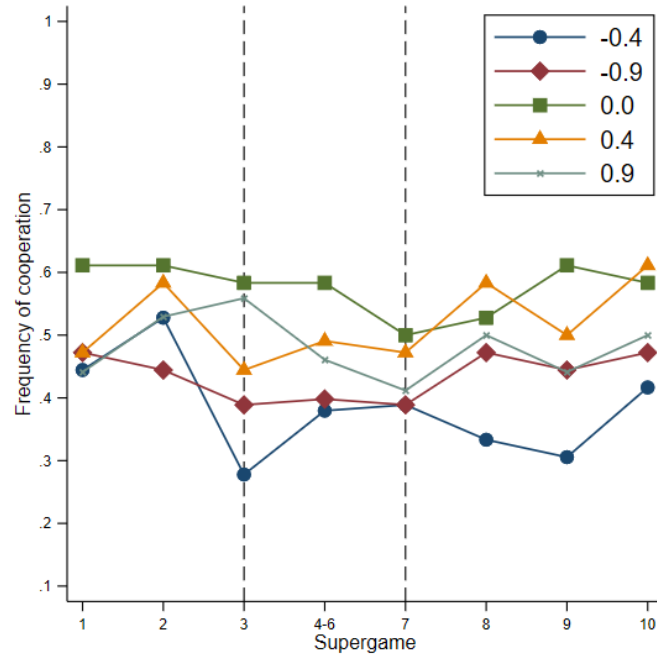
## 6 Results

To analyze the results, we first examine the link between the correlation structure and cooperation. This allows us to test Prediction 1 and, ultimately, offer an answer to Question 1 in Section 5. We then turn to the mechanisms through which correlation affects cooperation. First we examine the relationship between the structure of correlation and inferential error rates, which allows us to test Predictions 2-4 and answer Question 2 in Section 5. Finally, we test Prediction 5 regarding the link between the structure of correlation and strategies played by subjects, which allows us to address Question 3 in Section 5.

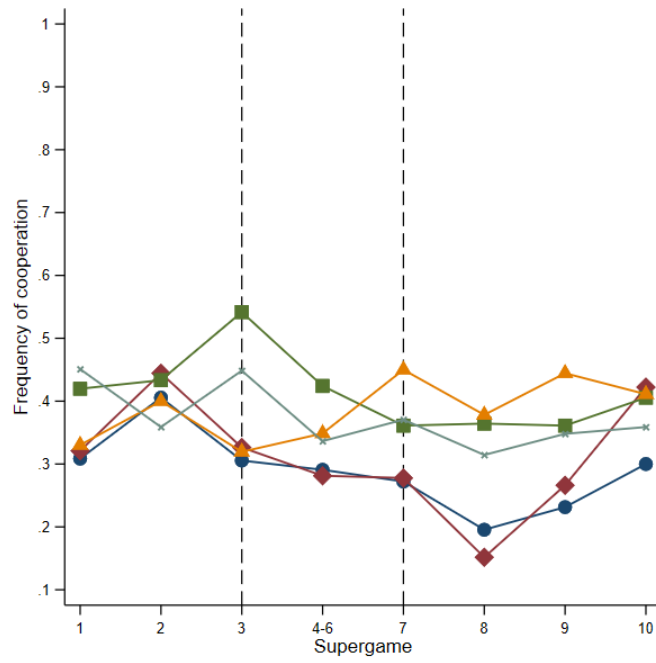
### 6.1 Correlation and Cooperation

Figure 2 shows the evolution of cooperation for all shock correlation structures: uncorrelated ( $\rho = 0$ ), strong negative correlation ( $\rho = -0.9$ ), moderate negative correlation ( $\rho = -0.4$ ), moderate positive correlation ( $\rho = 0.4$ ), and strong positive correlation ( $\rho = 0.9$ ). Average cooperation measures the proportion of rounds in which subjects cooperated in a supergame. We report results on cooperation for different correlation structures over subsequent supergames. We compare cooperation in the first stage of each supergame (Figure 2a) with cooperation in over all stages of those supergames (Figure 2b).

There is some level of cooperation for all correlation structures. A comparison between the top



(a) Round 1 Only



(b) All Rounds

Figure 2: Frequency of cooperation across supergames for all correlation structures



and bottom panels shows that cooperation seems higher in the first round than in subsequent rounds of a supergame, indicating that cooperation tends to unravel towards later rounds of the supergame. The evolution across supergames when all rounds are considered (Figure 2b) show a decline, albeit weak, in cooperation. Results in Figure 2 show that strong positive or negative correlation does not seem to induce more cooperation, relative to the baseline of no correlation. A pattern that contradicts our Prediction 1.

We report more formal measures of cooperation patterns across treatments (correlation structures) in Table 4. Once more, to better understand the evolution of cooperation, we examine cooperation in supergame 1, early stage supergames (1-3) and late stage supergames (7-10). For each block of supergames, we further disaggregate the results by first round and all rounds. Results in Table 4 confirm that, on average, cooperation is lower under strong correlation (both positive and negative) than when shocks are uncorrelated. They also confirm that cooperation generally unravels across rounds within a supergame and over subsequent supergames. The statistical significance reported in Table 4 refers to whether cooperation is statistically significantly different from zero.<sup>6</sup>

Table 4: Average Cooperation

| Rounds: | Supergame 1         |                     | Supergame 1-3       |                     | Supergame 7 - 10   |                     |
|---------|---------------------|---------------------|---------------------|---------------------|--------------------|---------------------|
|         | First               | All                 | First               | All                 | First              | All                 |
| -0.9    | 0.472**<br>(0.056)  | 0.321***<br>(0.003) | 0.435**<br>(0.061)  | 0.356***<br>(0.02)  | 0.444**<br>(0.062) | 0.222**<br>(0.043)  |
| -0.4    | 0.444***<br>(0.028) | 0.309***<br>(0.027) | 0.417**<br>(0.085)  | 0.335**<br>(0.05)   | 0.361<br>(0.135)   | 0.223**<br>(0.046)  |
| 0       | 0.611**<br>(0.073)  | 0.420**<br>(0.064)  | 0.602***<br>(0.037) | 0.451***<br>(0.040) | 0.556**<br>(0.092) | 0.367**<br>(0.057)  |
| 0.4     | 0.472***<br>(0.028) | 0.330*<br>(0.086)   | 0.500***<br>(0.042) | 0.347**<br>(0.038)  | 0.542**<br>(0.064) | 0.406***<br>(0.019) |
| 0.9     | 0.441***<br>(0.031) | 0.451***<br>(0.023) | 0.510**<br>(0.068)  | 0.425***<br>(0.036) | 0.463*<br>(0.143)  | 0.333<br>(0.156)    |

Notes: Robust standard errors (in parenthesis) are clustered at the session level. \* Indicates statistical significance at the 10% level (0.05  $\leq$  p-value  $\leq$  0.1). \*\* Indicates statistical significance at the 5% level (0.01  $\leq$  p-value  $\leq$  0.05). \*\*\* Indicates statistical significance at the 1% level (p-value  $\leq$  0.01)

Results in Table 4 do not offer information on the difference in cooperation rates across correlation structures. We present such information in Table 5. Table 5 shows the difference between cooperation with correlation and cooperation without correlation, as well as whether such differences

6. Unless otherwise stated, statistical significant is established by using a probit regression where errors are clustered at the session level.

are statistically significant. We report differences in cooperation rates for all rounds. Results in Table 5 show that negative correlation hinders cooperation, while positive correlation neither fosters nor hinders cooperation. The results are in stark contrast to Prediction 1.

Table 5: Difference in Average Cooperation (All Rounds)

|      | Supergame 1        | Supergame 1-3       | Supergame 7 - 10    |
|------|--------------------|---------------------|---------------------|
| -0.9 | -0.099*<br>(0.054) | -0.094**<br>(0.038) | -0.146**<br>(0.060) |
| -0.4 | -0.111*<br>(0.059) | -0.116**<br>(0.055) | -0.144**<br>(0.062) |
| 0.4  | -0.089<br>(0.091)  | -0.103**<br>(0.047) | 0.038<br>(0.051)    |
| 0.9  | 0.031<br>(0.058)   | -0.026<br>(0.046)   | -0.034<br>(0.1403)  |

Notes: This table gives the difference in average cooperation for each treatment from no correlation. Robust standard errors are in parentheses and are clustered at the session level

Given that our results indicate some learning over subsequent supergames, we focus on results from the last four supergames. These figures yield the main result of this paper, which answers Question 1 in Section 5:

**Result 1:** *On average, correlation does not improve cooperation.*

The average results conceal substantial heterogeneity across supergames. As discussed by Ostrom in her review of the literature (Ostrom 2000) a key force influencing cooperation in social dilemmas is the type of players involved in them and, specifically, the players’ willingness to engage in reciprocity that would lead to conditional cooperation. To investigate this in the context of our experiment, we tagged players in a way that is indicative of their willingness to initiate cooperation: whether they cooperated in the first round of each supergame or not. If a player cooperated in round 1, we labeled them as “nice” and they were labeled as “not nice” otherwise.<sup>7</sup>

Using our classification, we calculated the cooperation rate across supergames 7 – 10 for interactions between: (1) two nice players (‘nice-nice’), (2) two not nice players (‘not-not’), and (3) a nice player and a not nice player (‘nice-not’). Results are reported in Table 6. From these results we can see that when players are nice, positive correlation fosters cooperation from the baseline of  $\rho = 0$  ( $\rho = 0.9$   $p$ -value = 0;  $\rho = 0.4$   $p$ -value > 0.1). This is consistent with Prediction 1. This is not true when at least one not-nice player is present. And when both players are not-nice, correlation hinders cooperation ( $\rho = -0.9$   $p$ -value > 0.1;  $\rho = -0.4$   $p$ -value > 0.05).

7. We borrow this language from Bendor, Kramer, and Stout (1991), but our definition differs slightly from his. Bendor classified subjects as not-nice if they defected without a cause, i.e., without believing that their partner defected. However, we define not-nice as defecting in the very first round.

Table 6: Cooperation by Niceness for Supergames 7 - 10

| Treatment | Nice-Nice                      | Not-Not                       | Nice-Not                     |
|-----------|--------------------------------|-------------------------------|------------------------------|
| -0.9      | 0.581**<br>(0.097)<br>N = 308  | 0.079*<br>(0.020)<br>N=458    | 0.175<br>(0.069)<br>N=962    |
| -0.4      | 0.497*<br>(0.066)<br>N = 296   | 0.154***<br>(0.005)<br>N=846  | 0.186**<br>(0.020)<br>N=586  |
| 0         | 0.696**<br>(0.112)<br>N = 504  | 0.095*<br>(0.028)<br>N=326    | 0.282***<br>(0.027)<br>N=898 |
| 0.4       | 0.726**<br>(0.074)<br>N = 752  | 0.145<br>(0.068)<br>N=566     | 0.178*<br>(0.042)<br>N=410   |
| 0.9       | 0.955***<br>(0.009)<br>N = 332 | 0.024**<br>(0.005)<br>N = 410 | 0.244<br>(0.112)<br>N=890    |

Notes: Robust standard errors are in parentheses and are clustered at the session level

The results reported in Table 6 warrant a more qualified characterization of the relationship between correlation across shocks and cooperation. We provide such qualification in the following statement:

**Result 1’:** *When agents are nice (i.e., when agents are prone to cooperation in the initial rounds of the supergame), positive correlation fosters cooperation while negative correlation does not. When agents are not nice, strong correlation (either positive or negative), hinders cooperation.*

We now turn to the mechanisms underlying the link between correlation and cooperation to further elucidate the reasons why our Prediction 1 only seems to hold when players are nice.

## 6.2 Correlation and Inferential Error

In the previous section we reported evidence that, on average, higher levels of correlation does not foster higher levels of cooperation. We now further investigate the channel through which we anticipated the correlation structure influencing cooperation rates.

We used beliefs elicited from subjects using the BSR to estimate how inference error evolves from early supergames to later supergames. In our framework, after each round, subjects selected the probability with which they believed that the other player had cooperated. If they indicated a probability

greater than 0.5, we assigned their belief to cooperation (we assume they inferred cooperation). For probabilities less than 0.5, we assigned their belief to defection. We interpret a probability of 0.5 as the subject giving equal weighting to the probability to cooperation and defection. As such, we randomly assigned the subject’s belief to cooperation or defection with a 0.5 probability.<sup>8</sup> An inferential error occurs when there is a mismatch between the subject’s belief about their opponent’s action and the actual action taken by their opponent.

Table 7 shows the inferential error rates across correlation structures. Given that subjects do make errors outside of the region of uncertainty, we present the error rate across all actions (“All errors”) and the error rate when the player’s payoff falls within the region of uncertainty (“ROU errors”). We further disaggregated “All error” into Type 1 (incorrectly inferring defection) and Type 2 error (incorrectly inferring cooperation), and across early supergames (1-3) and late supergames (7-10).

As indicated in Table 7, subjects incur inferential errors across treatments and supergames (virtually all error rates are statistically significantly different from zero). They also incur errors when their payoff falls outside of the region of uncertainty. Given that subjects did incur errors outside of the region of uncertainty, albeit small relative to the region of uncertainty on average (8% outside versus 32% within the region of uncertainty), we focus on the total inferential error; that is, both within and outside of the region of uncertainty. Overall, across all treatments, inferential error is statistically greater than zero. However, inferential error is lower in late supergames than in early ones, which suggests learning by the subjects. Consequently, we base our discussion on results from later stage supergames where subjects have gained some experience.

Inferential error was smallest under  $\rho = 0.9$  (column 5 of Table 7), which is consistent with Prediction 2. The link between inferential error and negative correlation is non-monotonic. Moderate negative correlation raises inferential error (second row of column 5) and strong negative correlation reduces it (first row of column 5) relative to the baseline of no correlation. This is also consistent with Prediction 3. Therefore, strong correlation (either positive or negative, but especially the former) lowers inferential error.

To examine whether these differences in inferential error rates across treatments are distinguishable from chance, we calculated the difference between inferential error under each treatment (positive and negative correlation) and inferential error under the baseline (no correlation), as well as the standard deviation and statistical significance of these differences. Results are reported in Table 8. The results indicate that the effect of (strong) positive correlation on inferential error is distinguishable from chance, while the effect of negative correlation is not.

Negative correlation is less effective in strengthening monitoring (that is, reducing inferential error) than positive correlation because it expands the opponent’s domain of possible payoffs conditional on their actions, making it harder for a player to infer those actions (see distributions of player 2’s realized payoffs in Appendix C). To illustrate this, consider the following situation. If two players’ payoffs are positively correlated, and player 1 received a low payoff while player 2 received a high payoff, player 1 will have more confidence that the most likely action of player 2 was defection than if these payoffs were negatively correlated. If the payoffs are negatively correlated, this same situation

---

8. There were 439 such observations (out of a total of 14,952). After the random assignment, 50.57% of these were assigned as defection. Also, 27.3% of the observations fall within the region of uncertainty.

Table 7: Inferential Error Rates

| Treatment | Supergame 1-3      |                    |                    |                    | Supergame 7-10     |                    |                   |                   |
|-----------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|-------------------|-------------------|
|           | (1)                | (2)                | (3)                | (4)                | (5)                | (6)                | (7)               | (8)               |
|           | All Errors         | ROU Errors         | All Type 1         | All Type 2         | All Errors         | ROU Errors         | All Type 1        | All Type 2        |
| -0.9      | 0.20**<br>(0.03)   | 0.41***<br>(0.006) | 0.08**<br>(0.009)  | 0.12**<br>(0.023)  | 0.13*<br>(0.032)   | 0.28*<br>(0.067)   | 0.04*<br>(0.010)  | 0.09*<br>(0.023)  |
| -0.4      | 0.22***<br>(0.004) | 0.44***<br>(0.025) | 0.07***<br>(0.003) | 0.15***<br>(0.007) | 0.16**<br>(0.030)  | 0.34**<br>(0.044)  | 0.05**<br>(0.010) | 0.11**<br>(0.025) |
| 0         | 0.17**<br>(0.039)  | 0.38**<br>(0.061)  | 0.097<br>(0.039)   | 0.07***<br>(0.005) | 0.15**<br>(0.029)  | 0.32**<br>(0.054)  | 0.04**<br>(0.005) | 0.11*<br>(0.026)  |
| 0.4       | 0.20***<br>(0.011) | 0.36**<br>(0.038)  | 0.07**<br>(0.012)  | 0.12**<br>(0.020)  | 0.15***<br>(0.015) | 0.33***<br>(0.019) | 0.06**<br>(0.008) | 0.10**<br>(0.021) |
| 0.9       | 0.16*<br>(0.051)   | 0.30**<br>(0.036)  | 0.07*<br>(0.020)   | 0.08<br>(0.035)    | 0.06**<br>(0.010)  | 0.17**<br>(0.025)  | 0.02**<br>(0.004) | 0.04**<br>(0.007) |

Notes: Robust standard errors are in parentheses and are clustered at the session level

Table 8: Differences in Total Inferential Error (Supergames 7-10)

| Difference in error (SD) |                      |
|--------------------------|----------------------|
| -0.9                     | -0.020<br>(0.037)    |
| -0.4                     | 0.015<br>(0.037)     |
| 0.4                      | 0.004<br>(0.278)     |
| 0.9                      | -0.084***<br>(0.026) |

Notes: This table gives the difference in total error rates for each treatment relative to no correlation. Robust standard errors are in parentheses and are clustered at the session level

could arise for a variety of reasons. It could be that player 1 received the sucker payoff. Or it could be that player 2 had in fact cooperated but received a positive shock, while player 1 received a negative shock – a likely combination given that shocks are negatively correlated. This makes inference harder and raises the likelihood of inferential error. This is encapsulated in the following statement:

**Result 2:** *Stronger positive correlation reduces inferential error.*

The strength of monitoring itself, as measured by total inferential error, matters for cooperation. But this is not the only dimension of inferential error that matters. The composition of total inferential error is also crucial. As discussed in the development of Prediction 4, conditional on a player’s action, payoff received, and a correlation structure, there are two payoff distributions for her opponent: one when the opponent cooperates, and one when she defects. If these distributions overlap, the location of the benchmark within the overlapping region will determine the prevalence of Type I and II errors. This matters because if the benchmark is set at a level where most errors are Type I (Type II) errors, then the player will be less (more) likely to cooperate, perhaps inducing her opponent to defect (cooperate) more often. A high prevalence of Type I (Type II) errors will hinder (foster) cooperation. Given its importance for the level and composition of inferential errors, we now examine the extent to which players use the inferential process outlined in Prediction 4.

When a player cooperates, the center of the overlapping region of her opponent’s payoff distributions is around 54. And when she defects, it is around 19. Therefore, in line with Prediction 4, we expect that those are the payoffs the player will choose as benchmark values. While roughly in line with our Prediction (especially when considering variation across players), players tend to choose a lower value (on average) when they cooperate (46 with a standard deviation of 14) and a higher value (on average) when they defect (25 with a standard deviation of 12.5).

As also stated in Prediction 4, under the Bayesian updating rule outlined in Section 3, with positive (negative) values of  $\rho$ , subjects will more likely assume that their opponent defected (cooperated)

if they are signaled that their opponent’s payoff is above their benchmark value. For negative  $\rho$ s, the opposite applies. We find little evidence that subjects are updating their beliefs based on their private noisy signal. This rule explains subjects’ actions around 50% of the time, on average, across all treatments. The performance of this rule does not improve between supergames 1-3 and 7-10.

These results suggest a very limited role for the behavioral channel (agents use the signal to refine their inference) in the link between correlation and inferential error. Consequently, our results indicate that the reduction in inferential errors under stronger correlation is mostly driven by the mechanical channel, whereby correlation reduces inferential error without the subjects using the signal to update their beliefs. We summarize this in the following statement:

**Result 2:** *Subjects appear not to use their private signal to update their beliefs.*

Having found little evidence in favor of our predicted Bayesian updating process, we explore alternative ways in which subjects could be updating their beliefs. We explored a simple rule where subjects are more likely to assume cooperation if their own payoff is high enough and defection otherwise. Table 9 shows the average payoff within the region of uncertainty when subjects assume that the other player is defecting or cooperating.

Table 9: Average Payoff when Assuming Cooperation or Defection

|        |   | Assumption        |                   |
|--------|---|-------------------|-------------------|
|        |   | C                 | D                 |
| Action | C | 31.2/31<br>(3.56) | 42.9/43<br>(3.72) |
|        | D | 29.6/29<br>(3.68) | 41.9/42<br>(3.76) |

Notes: This table shows the mean/median payoff of player when they assume cooperation or defection of the other player. Standard deviation is in parenthesis.

The payoffs in Table 9 approximately coincide with the midpoint of the region of uncertainty of the player making the inference, as opposed to the center of the overlapping region of the opponent’s payoff distributions. The midpoints are 30.5 when cooperation is played and 42.5 when defection is played. This suggests the following inference rule. A subject will most likely assume cooperation if their payoff is substantially (more than 1 standard deviation) above the midpoint of their own region of uncertainty. Alternatively, a subject will most likely assume defection if their payoff is substantially (more than 1 standard deviation) below the midpoint of their own region of uncertainty.

However, if the player’s payoff is within 1 standard deviation (above or below) of the midpoint of their own region of uncertainty, they will most likely make the same assumption about the other player’s action as they did in the previous period. That is, there is a grey area where players will make inference based on how the history of play (or propensity to give others the benefit of the doubt if there is no history of play) shaped their perception of the other player’s type (either good or bad). We examine the extent to which this inferential rule (as opposed to the one advanced in Prediction 4) can rationalize inference drawn by subjects in the laboratory. Results of our analysis are reported in Table

10. This simple decision rule matches the data, on average, 88% of the time. For supergames 1-3, the match is 85%, and this increases to 90% in supergames 7-10. The match is high across treatments, which lends credence to this inferential process.

Table 10: Percentage Match of Simple Updating Rule

| All  |       |       | Supergame 1-3 |       | Supergame 7-10 |       |
|------|-------|-------|---------------|-------|----------------|-------|
|      | Mean  | SD    | Mean          | SD    | Mean           | SD    |
| -0.9 | 0.885 | 0.319 | 0.858         | 0.349 | 0.903          | 0.296 |
| -0.4 | 0.865 | 0.342 | 0.824         | 0.381 | 0.881          | 0.324 |
| 0    | 0.891 | 0.311 | 0.886         | 0.318 | 0.890          | 0.313 |
| 0.4  | 0.867 | 0.340 | 0.830         | 0.376 | 0.891          | 0.311 |
| 0.9  | 0.908 | 0.290 | 0.859         | 0.348 | 0.927          | 0.260 |

### 6.3 Correlation and Strategies Used by Subjects

In this section, we examine the relationship between correlated shocks and the nature of strategies used by players. Note that this mechanism can also be distinct from the effect of correlation on the inferential process. This is because a change in the correlation structure may prompt players to follow a more (or less) lenient strategy, conditional on whatever they infer about their opponent’s actions.

In Section 5, we predicted that a stronger positive correlation would raise the confidence with which players make inference and make them prone to employ more lenient strategies (Prediction 5). To elicit the most likely strategies played by subjects in the laboratory, we used the Strategy Frequency Estimation Method (SFEM) of Dal Bó and Fréchette (2011). In other words, we use a maximum likelihood estimation (MLE) approach to calculate the frequency of each strategy under consideration along with  $\beta$  that gives the model fit (see Appendix G for a detailed discussion). For this, we considered a subset of eight strategies from the twenty outlined in Fudenberg, Rand, and Dreber (2012). In Appendix H, we give a description of each. The strategies we consider consist of the eight top strategies identified by SFEM from data in Dal Bó and Fréchette (2018) (see Gill and Rosokha (2020)). The description of each strategy is similar to Fudenberg, Rand, and Dreber (2012).

We classified strategies into three categories: lenient, provocable and unfriendly. A lenient strategy (TF2T, AC, GRIM2) starts with cooperation and does not defect immediately when the opponent first defects. Provocable strategies (GRIM, 2TFT, TFT) start with cooperation but immediately defect when the opponent first defects. Unfriendly strategies (AD, DTFT) defect in the first round. The MLE estimation identifies that most prominent strategies and results are reported in Table 11. We focus our analysis on supergames 8-10, where subjects would have gained experience. Subjects predominantly used simple, memory-1 strategies. With memory-1 strategies, subjects only respond to their opponent’s action from the very last round. Memory-1 strategies include AD, DTFT, GRIM, TFT and AC.

Across all treatments, and on average, AD was the most frequent strategy, accounting for about half of the strategies used. Note that, the sum of the two most prominent provocable strategies, TFT



Table 11: Estimation of Strategies Used

| Correlation | Unfriendly         |                 | Provocable         |                  |                   | Lenient          |                    |                 |
|-------------|--------------------|-----------------|--------------------|------------------|-------------------|------------------|--------------------|-----------------|
|             | AD                 | DTFT            | GRIM               | 2TFT             | TFT               | GRIM2            | TF2T               | AC              |
| 0           | 0.38***<br>(0.119) | 0.05<br>(0.042) | 0.25***<br>(0.101) |                  | 0.16**<br>(0.088) | 0.06<br>(0.07)   | 0.07<br>(0.075)    | 0.03<br>(0.05)  |
| -0.4        | 0.57***<br>(0.131) | 0.09<br>(0.079) | 0.10*<br>(0.068)   |                  | 0.10<br>(0.084)   | 0.13*<br>(0.087) | 0.02<br>(0.028)    |                 |
| -0.9        | 0.52***<br>(0.087) | 0.02<br>(0.029) |                    | 0.12*<br>(0.074) | 0.15*<br>(0.095)  | 0.12*<br>(0.082) | 0.05<br>(0.06)     | 0.03<br>(0.045) |
| 0.4         | 0.38***<br>(0.099) | 0.03<br>(0.048) | 0.19**<br>(0.086)  | 0.02<br>(0.051)  | 0.08<br>(0.072)   | 0.06<br>(0.098)  | 0.24***<br>(0.093) |                 |
| 0.9         | 0.47***<br>(0.147) | 0.04<br>(0.061) | 0.11**<br>(0.061)  |                  | 0.15**<br>(0.081) | 0.08<br>(0.094)  | 0.08<br>(0.111)    | 0.07<br>(0.057) |
|             |                    |                 |                    |                  |                   |                  |                    |                 |

Notes: Bootstrapped standard errors are in parenthesis. Strategies that are 0.0 are dropped. \*\*\* p<0.01, \*\* p<0.05, \* p<0.1

and GRIM, is actually higher than AD when  $\rho = 0$ . Yet, AD becomes more prominent (even more than the sum of provokable strategies) when correlation got stronger. This rejects our Prediction 5 that subjects will play more lenient strategies as positive correlation becomes stronger. This result, however, is in line with the fact that cooperation was, on average, lower under stronger correlation. Also, in line with cooperation patterns, AD is more prominent under negative correlation. These results are summarized in the following insight.

**Result 4:** *Subjects played mostly unfriendly strategies, regardless of the correlation structure.*

Results 1-4 offer mixed evidence regarding our predictions. Predictions 2 and 3 are supported by the data, while Prediction 4 is not. And Prediction 1 is rejected, but a qualified version of it is supported by the data (Result 1'). In the next section, we combine these pieces of evidence to better elucidate the nature of the relationship between correlation and cooperation.

## 6.4 Mechanisms and the Effect of Correlation on Cooperation

A crucial aspect of our results is that a strong positive correlation leads to enhanced monitoring on average, as revealed by lower inferential errors (Table 8). And this does not translate into more cooperation, on average (Table 5). However, it decidedly fosters cooperation when both agents are nice (that is, when they cooperate in the early stages of the game) and hinders it when both agents are not nice (Table 6). One possible explanation for this result is that, once we disaggregate by type of players, strong positive correlation lowers inferential error only when players are nice, thereby strengthening monitoring and, consequently, cooperation. Another possible explanation is that enhanced monitoring only fosters cooperation if players are nice.

To explore the first possible explanation, we examine inferential error under alternative correlation structures; but this time we disaggregate results by the type of players involved in the game. The results are reported in Table 12. When we compare our results with those for the average situation (Table 8), we can see that virtually the same patterns emerge – strong (but not moderate) positive correlation lowers inferential errors when players are both nice and when they are both not-nice. It also lowers inferential error but to a lesser extent when one player is nice and the other is not. In contrast, negative correlation (either moderate or strong) does not lower inferential error in a way that is distinguishable from chance. This is with the exception of high negative correlation when both players are not nice.

Results in Table 12, leaves us with the second possible explanation – that enhanced monitoring only fosters cooperation if players are nice. The results seem to strongly support this explanation. This follows from the observation that the reduction in inferential error is largest when correlation is strong and positive. In fact, inferential errors when correlation is 0.9 is almost indistinguishable from zero. Coincidentally, cooperation is also most prevalent (in fact, highly prevalent at 95%) when correlation is 0.9. This strongly suggests that players' ability to monitor their opponent (and players' knowledge about their opponents' ability to monitor them) disciplines the subjects into cooperation. But this insight is much more nuanced than may seem at first glance. As correlation raises from 0 to 0.9, cooperation decidedly drops from 10% to 2.5% (Table 6), despite inferential error dropping from 24% to 6% (Table 12).

Table 12: Error Rates by Niceness for Supergames 7 – 10

| Treatment    | Nice-Nice          | Not-Not             | Nice-Not            |
|--------------|--------------------|---------------------|---------------------|
| -0.9         | 0.107**<br>(0.015) | 0.107<br>(0.05)     | 0.146**<br>(0.032)  |
| -0.4         | 0.118**<br>(0.005) | 0.199***<br>(0.012) | 0.135*<br>(0.032)   |
| 0            | 0.113*<br>(0.030)  | 0.239**<br>(0.04)   | 0.135***<br>(0.012) |
| 0.4          | 0.122*<br>(0.031)  | 0.208*<br>(0.056)   | 0.127***<br>(0.007) |
| 0.9          | 0.018<br>(0.012)   | 0.059**<br>(0.008)  | 0.084**<br>(0.017)  |
| Observations | N = 332            | N = 410             | N = 890             |

Robust standard errors are in parentheses and are clustered at the session level

Our results suggest that a strong positive correlation improves monitoring, like we predicted. But in contrast to our prediction, improved monitoring does not necessarily foster cooperation – it simply better reveals to players the actions of their opponents. If those actions happen to be non-cooperative, as it tends to be the case at first if their opponent is not nice, then they move more quickly towards defection to avoid being the sucker. On the other hand, if the actions revealed by improved monitoring happen to be cooperative, as it tends to be the case at first if their opponent is nice, then they move more decisively towards cooperation.

In summary, our results indicate that a strong positive correlation removes the veil of ignorance regarding the actions of the opponent; but its effect on cooperation depends on whether that reveals a cooperative or non-cooperative opponent. If stronger monitoring reveals an opponent that is reluctant to cooperate, then it will prompt cooperation to unravel. Of course, a player should know that improved monitoring will make her actions more visible as well, and that consistent defection will induce her opponent to retaliate, thereby condemning her to a low payoff. So, why are “not nice” players not anticipating this? If the player thinks her opponent will not give her the opportunity to build trust and cooperation, she may swiftly move to defect in anticipation of a bad equilibrium. As long as “not nice” players are in the mix and they are revealed by improved monitoring, there is a chance that players will defect to avoid a sucker payoff.

## 7 The Computational Experiment

In Section 6.2. we conjectured that, as long as not nice players are in the mix, improved monitoring (from strong positive correlation) may lower cooperation. This is because players may be pessimistic about the prospects of cooperation when improved monitoring reveals a not nice opponent, and defect preemptively to avoid becoming the sucker in the prisoner’s dilemma game. We are not the first ones to identify the composition (that is, types) of players involved in a repeated game as an important force underlying the prevalence of cooperation. The presence of not nice players may increase the prevalence of a bad (defection) equilibrium. Ostrom (2000) made this point in her review of the literature on cooperation in social dilemmas. But she, along with Axelrod (1980), also surmised that repeated interactions may stimulate an evolutionary process by which not nice players are selected out of the pool. To investigate this further, we complement our laboratory experiment with a computational simulation of such evolutionary process.

We use a genetic algorithm (Holland 1975) for the evolutionary process. The genetic algorithm starts with a pool of candidate strategies, called automata. The interaction among automata mirrors the experimental framework closely, but there are some notable differences. First, instead of playing ten supergames, automata interact over hundreds of generations. In one generation each automaton is matched to every other automaton, including itself, and they play a supergame. Second, from one generation to the next, the environment dynamically changes given that only strategies with the best performance survive (that is, advance to the next generation).

Strategies also undergo a process of mutation. This process is not removed from reality. One can consider this as a process of exploration and learning, where the agents interacting in the environment try new strategies. If successful, then these strategies are mimicked by other players. If a strategy is unsuccessful, it is abandoned and no longer part of the strategy pool in subsequent generations. Lastly, the automata update their beliefs on the likely action of their opponent according to the decision rule we characterized in Section 3. This is another way in which automata differ from subjects in the laboratory since, from the experimental results, we found that subjects did not seem to follow this rule closely.

Despite these differences, we believe that this is an appropriate framework to gain insights into how cooperation can evolve over many repeated interactions under various correlation structures. Genetic algorithms have been widely used to show how strategies evolve under various environments. We use the tournament style pioneered by (Axelrod 1980) that have since then been used to examine the evolution of strategies in repeated noisy PDs. Similar computational exercises, by and large, have focused on implementation and perception errors (Ioannou 2014b; Miller 1996; Zhang 2018). Other repeated PD applications include examining the impact of costly strategy adjustments (Romero and Rosokha 2019).

### 7.1 The Computational Experiment

Our evolutionary process is similar to Miller (1996) and Ioannou (2014b). Strategies are implemented as finite state automata, with a fixed number of states. Each internal state is accorded cooperation

(C) or defection (D). Conditioned on what the other player does, each state has transition rules that dictate what state the automaton should next transition to. Given that subjects mostly used memory-1 strategies, and to reduce computational time, we focus only on finite automata with at most two states. Memory-1 strategies only consider the immediate past move of its opponent. In our environment, each automaton carries information about not only the strategies, but also, the benchmark values relative to which the signal is defined. The two benchmark values that are needed if the automaton is in a cooperation or defection state are endogenously determined; that is, the benchmark values must survive the evolutionary process just like the strategy must. As such, each automaton is represented by a 21-bit binary string. The first 7-bit translates into an integer in the range [24, 84] and the second 7-bit an integer in the range [-11, 49]. The final 7-bits represent the strategy. For this last 7-bit, the first bit represents the starting state, and the rest gives the transitional states. While there are 128 representations of 7-bit strings, these map into 26 unique strategies. We provide more details on the representation of the automata in Appendix I.

For the evolutionary process (summarized in Table 13), we begin with a population of thirty randomly generated 21-bit binary strings. In every generation, an automaton is paired with every other automaton, including a copy of itself, to play the repeated PD for eighty rounds on average (we use a continuation probability of 0.9875). In the first generation, the automaton also randomly selects two benchmark values, one is used when it cooperates and the other when it defects. At the end, a performance score is calculated based on average payoff across all interactions. If there are ties, one is randomly selected. To populate the next generation, the top twenty performers are selected (the parents). These parents create ten children through a process of mutation. A pair is randomly selected from the parents, and the top performer of this pair undergoes mutation. Mutation involves changing every bit of each string with a probability of 0.04. Mutation occurs on both the strategy and the two benchmark values. Therefore, the automaton updates both its strategy and its benchmark values. This method of mutation does not guarantee that all the best of the best strategies will be selected, or that all the worst of the best will be eliminated. The evolutionary process is simulated for 1000 generations. Then the entire process is repeated 100 times.

Table 13: The Evolutionary Process of Genetic Algorithm

|         |  |
|---------|--|
| Step 1: | Initialize Generation T0 of 30 randomly generated automata   |
| Step 2: | Initiate round-robin tournament  |
| Step 3: | For $t = 0$ to $T$ :<br>For each automaton:<br>For round 1 to length of round:<br>Play all other automata and self<br>Determine average payoff using PD matrix |
| Step 4: | Select 20 automata with highest payoff   |
| Step 5: | Create 10 - pairs:<br>Select from each pair, automata with highest payoff  |
| Step 6: | $t = t + 1$ Mutate every bit with probability of 0.04  |

## 7.2 Results of Computational Experiment

Across all levels of correlation, cooperation begins at approximately 50% (Figure 3). This is expected given that the automata are randomly selected in the first generation. In fact, out of the 126 possible strategies, 40 are AC and 40 are AD. Immediately, cooperation suffers a sharp reduction. But as generations progress, cooperation rises and stabilizes at higher levels. Higher levels of correlation, both positive and negative, induce higher levels of cooperation. Initially, and for a large number of generations, cooperation under  $\rho = -0.9$  is lower than under  $\rho = 0.9$ . However, cooperation levels under  $\rho = -0.9$  and  $\rho = 0.9$  converge by generation 1000. Notably, all levels of  $\rho$  yield higher cooperation rates than  $\rho = 0$ , except for  $\rho = -0.4$ . These results are consistent with our Predictions 1-4. Correlation structures that lower inferential error (both moderate and strong positive correlation, as well as strong negative correlation), also foster cooperation, albeit in the long run, and as a result of the evolutionary process. Negative moderate correlation, which we both predicted and found evidence to support that it raises inferential error, hinders cooperation.

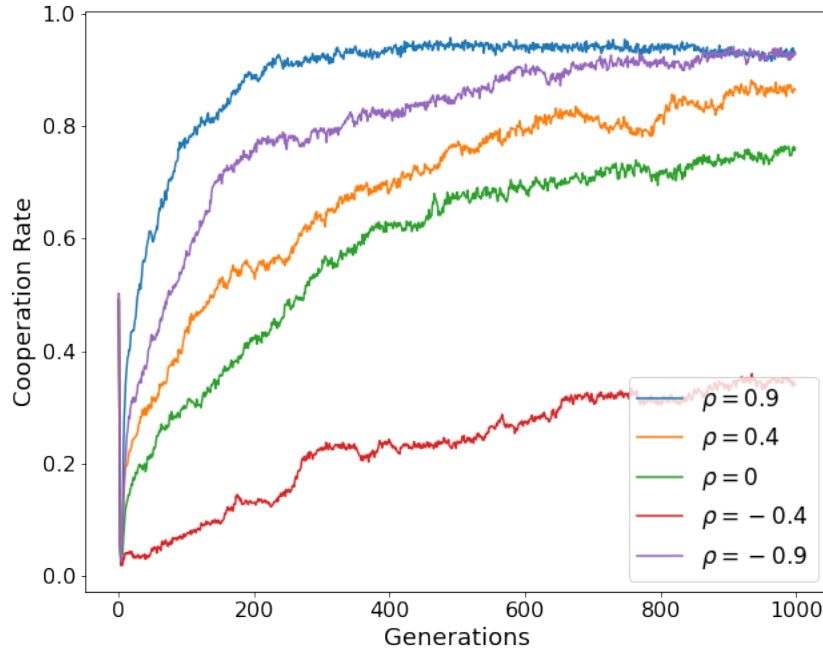


Figure 3: Evolution of cooperation over 1000 generations

As discussed before, the correlation structure may affect cooperation through the inferential error, and subsequently influence strategies employed by players. Figure 4 shows the evolution of strategies for each correlation structure. We display the memory-1 strategies that the MLE indicated best matched the subjects' data. Mostly two strategies, TFT and AD, explain the differences in cooperation rates observed. For  $\rho = 0.9$ , where the highest cooperation rate was observed, AD very quickly dies out in the population, leaving TFT as the most frequent strategy. The lower the

cooperation rate, the slower AD disappears from the population. This suggests that with sustained interaction, cooperation can be maintained. For high levels of  $\rho$ , agents will converge to high levels of cooperation quickly. This process is slower for moderate levels of correlation and  $\rho = 0$ .

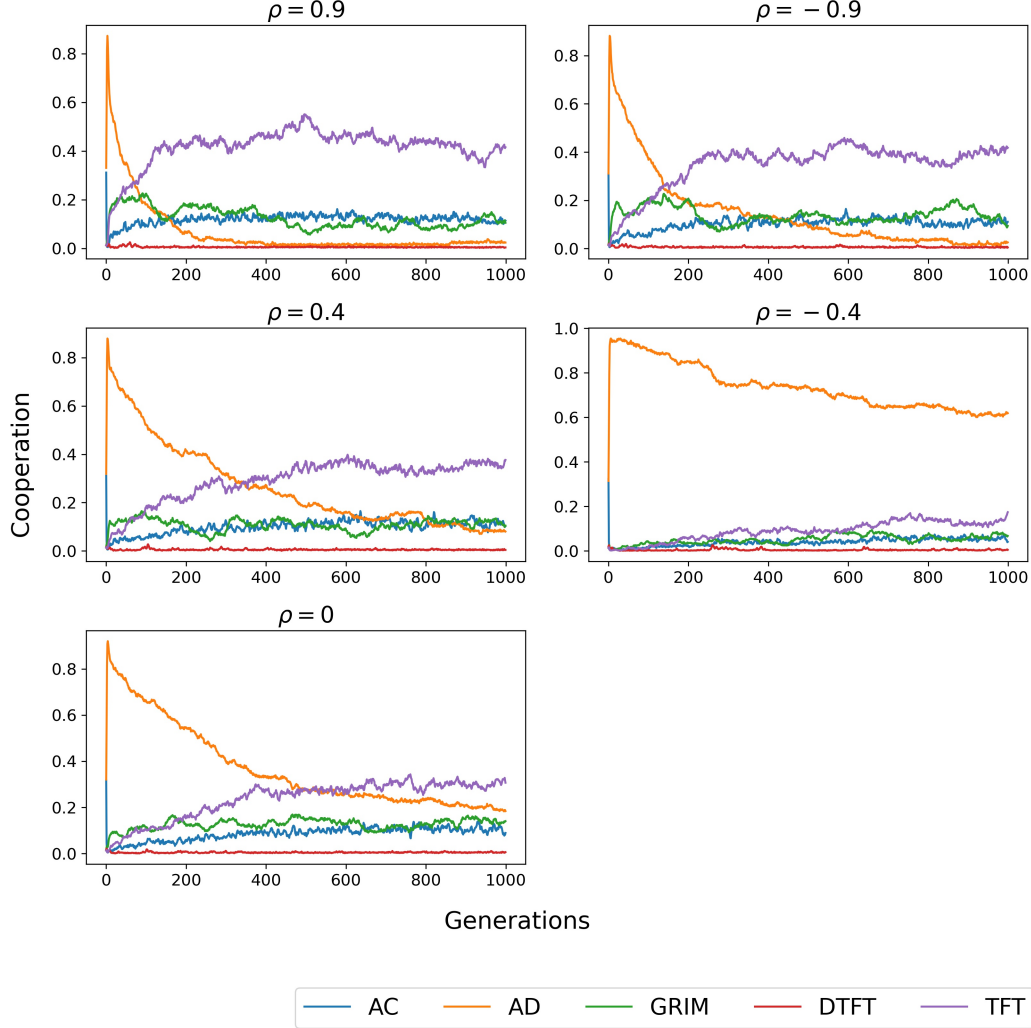


Figure 4: Evolution of strategies under each correlation

The evolution of the benchmarks offers additional insight into the effect of correlation on cooperation. We report those results in Figure 5. The benchmark value ascribed to cooperation under positive correlation and under the baseline increased to around 70 by generation 500. The benchmark value falls quickly to 35 under strong negative correlation, but only to 45 under moderate negative correlation. Recall that the decision rule prescribes cooperation if there is a signal that the other player's realized payoff is below (above) the benchmark value for  $\rho \geq 0$  ( $\rho < 0$ ). Therefore, a (low) high benchmark value under positive (negative) correlation implies, all else constant, that it is more likely the automaton will cooperate in error, rather than defect in error.

The pattern of the rate at which the average benchmark converges to its long-term value, matches

the cooperation rate patterns. Under a positive (negative) and strong correlation the benchmark converges more quickly and to a higher (lower) long-term value, and this coincides with a quicker convergence to cooperation (Figure 3).

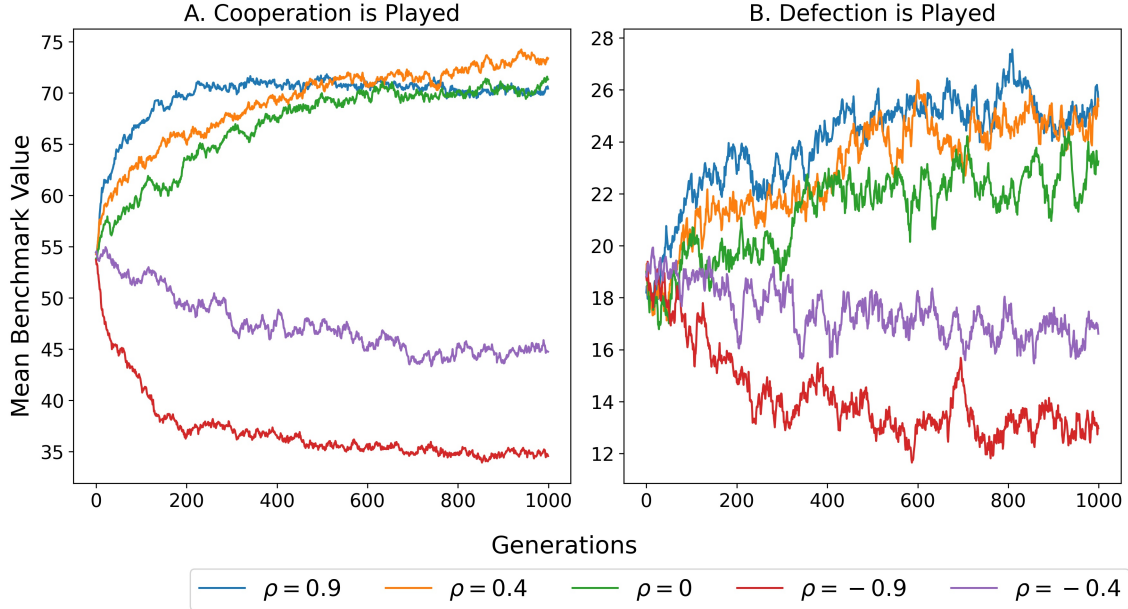


Figure 5: Evolution of benchmark values under each correlation

A closer look at results in Figures 4 and 5 points to an important interaction between both mechanisms underlying the effect of correlation on cooperation: the strategy mechanism and the inferential error mechanism. At the long-term benchmark values observed in Figure 5, conditional strategies become somewhat unresponsive to inferential error. When an agent’s realized payoff falls within the region of uncertainty, they are more likely to cooperate than defect under the long-term benchmark values. This is only reinforced when the dominant strategy is TFT. If the automaton employs a TFT strategy and selects a high (low) benchmark under positive (negative) correlation, it is more likely to cooperate in error, than to defect in error. This fosters cooperation even further and offers an additional explanation to the high level of cooperation in the computational experiment relative to the laboratory experiment.

## 8 Discussion

In this study, we examined the impact of correlation on cooperation in an infinitely repeated prisoner’s dilemma. In the stage game, players’ payoffs are affected by a random shock that is uniformly distributed. This shock is independent across rounds but correlated across players. We explored five correlation structures, two positive correlation levels (moderate and high), two negative correlation levels (moderate and high), and a baseline of no correlation. We found that, on average, correlation does not enhance cooperation. While negative correlation weakly lowers cooperation, positive correlation has no impact on cooperation.



We offered two explanations for this observation. First, we anticipated that higher levels of correlation would improve cooperation through two mechanisms: a reduction in inferential error when correlation is positive or negative and strong (including a purely mechanical channel and a behavioral channel), and a change in strategies towards more lenient alternatives. We found evidence supporting the first mechanism, albeit weaker than predicted because subjects in the lab failed to fully exploit the behavioral channel. Second, and also on average, more lenient strategies did not become more prevalent under stronger correlation.

At first glance, it seems puzzling that improved monitoring delivered by positive and strong negative correlation structures does not translate into more cooperation on average. However, the reason for this becomes apparent when we explore the heterogeneity concealed in the average effect. We found that improved monitoring associated with certain correlation structures may have simply revealed the uncooperative nature of many players in the subject pool. In sessions where players were prone to cooperate (defect) at the beginning of the supergame, positive and strong negative correlation structures led to lower inferential error and higher (lower) cooperation. This observation lends support to the argument that cooperation in social dilemmas depends greatly on the environment. In our case, we see that allowing for better monitoring does not automatically lead to cooperation when monitoring is imperfect.

Our results give preliminary insights into the prospects of cooperation in groups with correlated outcomes. Noisy payoffs that are uncorrelated introduce an environment of imperfect monitoring. This imperfect monitoring environment seems to create a veil of uncertainty that encourages free-riding. It appears that individuals are inclined to act in their self-interest if they believe that their uncooperative behavior can go unnoticed. As better monitoring partially removes this veil of uncertainty, cooperation strengthens (unravels) if players are prone (reluctant) to cooperate at initial stages of the game. As pointed out before, there are many environments where the correlation structure of shocks is important to effectiveness of economic groups – perhaps more prominently environments of mutual insurance. Our results suggest that risk-sharing and cooperation can interact in complicated ways depending on the composition of the subject pool.

We also found some evidence that negative correlation weakens cooperation in comparison to the baseline of no correlation. This is true even in the presence of an evolutionary process that tends to select uncooperative players out of the pool. This seems problematic for economic groups that provide mutual insurance. Negatively correlated shocks mean that whenever someone in the group has had a bad shock, another member has had a good shock to offset this. This sounds ideal for risk-sharing. However, we found that moderate negative correlation introduces additional noise, as players found it more difficult to unravel the information contained in the correlation structure, relative to positively correlated shocks. In these circumstances, there appears to be a trade-off between cooperation and risk-sharing. In light of this trade-off, it seems preferable to include agents with uncorrelated payoffs into the group, instead of agents with negatively correlated payoffs. Future work should explicitly consider the impact of correlated shocks in a risk-sharing environment.

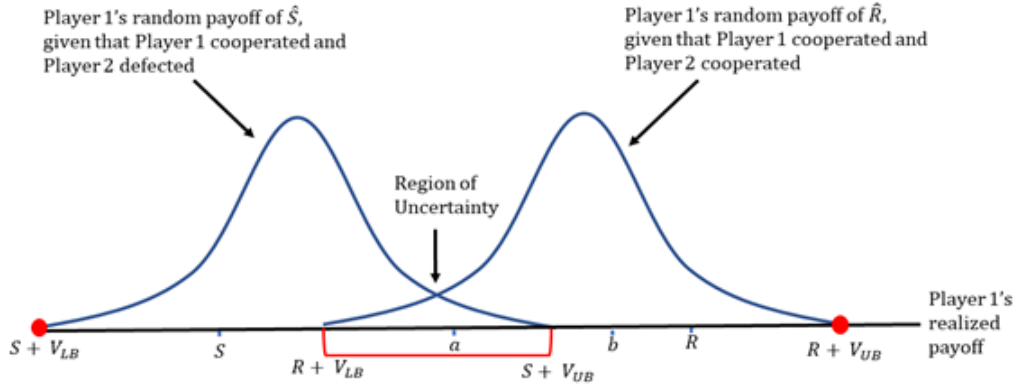
## References

- Agarwal, Bina. 2018. “Can group farms outperform individual family farms? Empirical insights from India.” *World Development* 108:57–73.
- Aoyagi, Masaki, V Bhaskar, and Guillaume R Fréchette. 2019. “The impact of monitoring in infinitely repeated games: Perfect, public, and private.” *American Economic Journal: Microeconomics* 11 (1): 1–43.
- Axelrod, Robert. 1980. “Effective choice in the prisoner’s dilemma.” *Journal of conflict resolution* 24 (1): 3–25.
- Baland, Jean-Marie, Isabelle Bonjean, Catherine Guirking, and Roberta Ziparo. 2016. “The economic consequences of mutual help in extended families.” *Journal of Development Economics* 123:38–56.
- Baland, Jean-Marie, Catherine Guirking, and Charlotte Mali. 2011. “Pretending to be poor: Borrowing to escape forced solidarity in Cameroon.” *Economic development and cultural change* 60 (1): 1–16.
- Bendor, Jonathan. 1993. “Uncertainty and the evolution of cooperation.” *Journal of Conflict resolution* 37 (4): 709–734.
- Bendor, Jonathan, Roderick M Kramer, and Suzanne Stout. 1991. “When in doubt... Cooperation in a noisy prisoner’s dilemma.” *Journal of conflict resolution* 35 (4): 691–719.
- Bloch, Francis, Garance Genicot, and Debraj Ray. 2008. “Informal insurance in social networks.” *Journal of Economic Theory* 143 (1): 36–58.
- Chen, Daniel L, Martin Schonger, and Chris Wickens. 2016. “oTree—An open-source platform for laboratory, online, and field experiments.” *Journal of Behavioral and Experimental Finance* 9:88–97.
- Dal Bó, Pedro, and Guillaume R Fréchette. 2011. “The evolution of cooperation in infinitely repeated games: Experimental evidence.” *American Economic Review* 101 (1): 411–29.
- . 2018. “On the determinants of cooperation in infinitely repeated games: A survey.” *Journal of Economic Literature* 56 (1): 60–114.
- . 2019. “Strategy choice in the infinitely repeated Prisoner’s Dilemma.” *American Economic Review* 109 (11): 3929–52.
- Danz, David, Lise Vesterlund, and Alistair J Wilson. 2020. *Belief elicitation: Limiting truth telling with information on incentives*. Technical report. National Bureau of Economic Research.
- Fafchamps, Marcel. 2011. “Risk sharing between households.” *Handbook of social economics* 1:1255–1279.
- Fitzsimons, Emla, Bansi Malde, and Marcos Vera-Hernández. 2018. “Group size and the efficiency of informal risk sharing.” *The Economic Journal* 128 (612): F575–F608.

- Fudenberg, Drew, David G Rand, and Anna Dreber. 2012. "Slow to anger and fast to forgive: Cooperation in an uncertain world." *American Economic Review* 102 (2): 720–49.
- Fundenberg, Drew, and Eric Maskin. 1990. "Evolution and cooperation in noisy repeated games." *The American Economic Review* 80 (2): 274–279.
- Gill, David, and Yaroslav Rosokha. 2020. "Beliefs, learning, and personality in the indefinitely repeated prisoner's dilemma." *Available at SSRN 3652318*.
- Holland, John H. 1975. *Adaptation in natural and artificial systems: An introductory analysis with applications to biology, control, and artificial intelligence*. Adaptation in natural and artificial systems: An introductory analysis with applications to biology, control, and artificial intelligence. Pages: viii, 183. Oxford, England: U Michigan Press. ISBN: 978-0-472-08460-9.
- Hossain, Tanjim, and Ryo Okui. 2013. "The binarized scoring rule." *Review of Economic Studies* 80 (3): 984–1001.
- Imhof, Lorens A, Drew Fudenberg, and Martin A Nowak. 2007. "Tit-for-tat or win-stay, lose-shift?" *Journal of theoretical biology* 247 (3): 574–580.
- Ioannou, Christos A. 2014a. "Asymptotic behavior of strategies in the repeated prisoner's dilemma game in the presence of errors." *Artif. Intell. Res.* 3 (4): 28–37.
- . 2014b. "Coevolution of finite automata with errors." *Journal of Evolutionary Economics* 24 (3): 541–571.
- Jakiela, Pamela, and Owen Ozier. 2016. "Does Africa need a rotten kin theorem? Experimental evidence from village economies." *The Review of Economic Studies* 83 (1): 231–268.
- Kayaba, Yutaka, Hitoshi Matsushima, and Tomohisa Toyama. 2020. "Accuracy and retaliation in repeated games with imperfect private monitoring: Experiments." *Games and Economic Behavior* 120:193–208.
- Miller, John H. 1996. "The coevolution of automata in the repeated prisoner's dilemma." *Journal of Economic Behavior & Organization* 29 (1): 87–112.
- Ostrom, Elinor. 2000. "Collective action and the evolution of social norms." *Journal of economic perspectives* 14 (3): 137–158.
- Ostrom, Elinor, Joanna Burger, Christopher B Field, Richard B Norgaard, and David Policansky. 1999. "Revisiting the commons: local lessons, global challenges." *science* 284 (5412): 278–282.
- Rand, David G, Drew Fudenberg, and Anna Dreber. 2015. "It's the thought that counts: The role of intentions in noisy repeated games." *Journal of Economic Behavior & Organization* 116:481–499.
- Romero, Julian, and Yaroslav Rosokha. 2019. "The evolution of cooperation: The role of costly strategy adjustments." *American Economic Journal: Microeconomics* 11 (1): 299–328.
- Zhang, Huanren. 2018. "Errors can increase cooperation in finite populations." *Games and Economic Behavior* 107:203–219.

## Appendix A The Player's Own Payoff Distributions

This is the distribution of realized payoff for player when she cooperates according to Bendor (1993). The shocks faced by each player is independent. The size of the region of uncertainty is determined by the variance of the shock faced by the player. In Bendor (1993), each player sets critical cutoff values such that Type 1 and Type 2 errors have the same probability  $p$  of occurring, where  $\frac{1}{2} > p > 0$ . Also, each player knows where the realized payoff of their opponent lies in relation to the critical cutoff value.



## Appendix B Bayesian Updating Process

|                 | Signal       |                     |
|-----------------|--------------|---------------------|
|                 | $s = 0$      | $s = 1$             |
| State of Nature | $\theta = C$ | $\pi_c$ $1 - \pi_c$ |
|                 | $\theta = D$ | $1 - \pi_D$ $\pi_D$ |

Let  $s = 0$  be that player 1 received a signal that player 2's realized payoff is above the benchmark value and  $s = 1$  be that it is below.

Using a Bayesian approach, player 1 will update her prior after she receives a signal according to:

$$P(\theta = C/s = 0) = \frac{P(\theta = C \cap s = 0)}{P(s = 0)}$$

$$P(\theta = C/s = 0) = \frac{P(s = 0/\theta = C)P(\theta = C)}{P(s = 0)}$$

Directly from Bayes' rule:

$$P(\theta = C/s = 0) = \frac{\pi_C P(\theta = C)}{\pi_C P(\theta = C) + (1 - \pi_C)(1 - P(\theta = C))}$$

Where  $P(\theta = C) = 0.5$ . This is player 1's prior belief on the probability of cooperation.  $\pi_C$  ( $P(s = 0/\theta = C)$ ) can be thought of as player 1's belief of the type of signal that is possible when the other player cooperates. With these information, player 1 can calculate  $P(\theta = C/s = 0)$ , the probability that the true state is cooperation given that there is signal that player 2 is above the benchmark value.

*The decision rule:*

A signal is informative, or support your belief about the state, if and only if  $P(\theta = C/s = 0) > P(\theta = D/s = 0)$  and  $P(\theta = C/s = 1) < P(\theta = D/s = 1)$ . That is  $\pi_C > 1 - \pi_D$  and  $1 - \pi_C < \pi_D$ . Using the former expression, a signal is informative about a state if it is more likely to occur in the cooperation state and not very likely to occur in the defection state.

If player 1 is cooperating and gets a signal that player 2 is above a benchmark value, Table 14 shows Player 1's beliefs about Player 2's most likely action.

| $\rho$ | $RP$ | $T1$ | $T2$ | $\pi_C$ | $1 - \pi_C$ | $\pi_D$ | $1 - \pi_D$ | Error | P2's most likely action |
|--------|------|------|------|---------|-------------|---------|-------------|-------|-------------------------|
| 0.9    | 24   | 0.00 | 0.03 | 0.00    | 1.00        | 0.03    | 0.97        | 0.03  | D                       |
| 0.4    | 24   | 0.06 | 0.25 | 0.06    | 0.94        | 0.25    | 0.75        | 0.31  | D                       |
| 0.0    | 24   | 0.36 | 0.35 | 0.36    | 0.64        | 0.35    | 0.65        | 0.71  | D                       |
| -0.4   | 24   | 0.23 | 0.55 | 0.77    | 0.23        | 0.45    | 0.55        | 0.78  | C                       |
| -0.9   | 24   | 0.00 | 0.34 | 1.00    | 0.00        | 0.66    | 0.34        | 0.34  | C                       |
| 0.9    | 25   | 0.00 | 0.02 | 0.00    | 1.00        | 0.02    | 0.98        | 0.02  | D                       |
| 0.4    | 25   | 0.10 | 0.25 | 0.10    | 0.90        | 0.25    | 0.75        | 0.35  | D                       |
| 0.0    | 25   | 0.36 | 0.36 | 0.36    | 0.64        | 0.36    | 0.64        | 0.72  | D                       |
| -0.4   | 25   | 0.31 | 0.55 | 0.69    | 0.31        | 0.45    | 0.55        | 0.86  | C                       |
| -0.9   | 25   | 0.00 | 0.26 | 1.00    | 0.00        | 0.74    | 0.26        | 0.26  | C                       |
| 0.9    | 26   | 0.00 | 0.01 | 0.00    | 1.00        | 0.01    | 0.99        | 0.01  | D                       |
| 0.4    | 26   | 0.13 | 0.24 | 0.13    | 0.87        | 0.24    | 0.76        | 0.37  | D                       |
| 0.0    | 26   | 0.36 | 0.36 | 0.36    | 0.64        | 0.36    | 0.64        | 0.72  | D                       |
| -0.4   | 26   | 0.36 | 0.52 | 0.64    | 0.36        | 0.48    | 0.52        | 0.88  | C                       |
| -0.9   | 26   | 0.00 | 0.21 | 1.00    | 0.00        | 0.79    | 0.21        | 0.21  | C                       |
| 0.9    | 27   | 0.00 | 0.01 | 0.00    | 1.00        | 0.01    | 0.99        | 0.01  | D                       |
| 0.4    | 27   | 0.14 | 0.23 | 0.14    | 0.86        | 0.23    | 0.77        | 0.37  | D                       |
| 0.0    | 27   | 0.36 | 0.36 | 0.36    | 0.64        | 0.36    | 0.64        | 0.72  | D                       |
| -0.4   | 27   | 0.40 | 0.51 | 0.60    | 0.40        | 0.49    | 0.51        | 0.91  | C                       |
| -0.9   | 27   | 0.01 | 0.18 | 0.99    | 0.01        | 0.82    | 0.18        | 0.19  | C                       |
| 0.9    | 28   | 0.00 | 0.00 | 0.00    | 1.00        | 0.00    | 1.00        | 0.00  | D                       |
| 0.4    | 28   | 0.16 | 0.21 | 0.16    | 0.84        | 0.21    | 0.79        | 0.37  | D                       |
| 0.0    | 28   | 0.35 | 0.36 | 0.35    | 0.65        | 0.36    | 0.64        | 0.71  | D                       |
| -0.4   | 28   | 0.41 | 0.50 | 0.59    | 0.41        | 0.50    | 0.50        | 0.91  | C                       |
| -0.9   | 28   | 0.02 | 0.16 | 0.98    | 0.02        | 0.84    | 0.16        | 0.18  | C                       |
| 0.9    | 29   | 0.00 | 0.00 | 0.00    | 1.00        | 0.00    | 1.00        | 0.00  | D                       |
| 0.4    | 29   | 0.17 | 0.21 | 0.17    | 0.83        | 0.21    | 0.79        | 0.38  | D                       |

|      |    |      |      |      |      |      |      |      |   |
|------|----|------|------|------|------|------|------|------|---|
| 0.0  | 29 | 0.36 | 0.36 | 0.36 | 0.64 | 0.36 | 0.64 | 0.72 | D |
| -0.4 | 29 | 0.44 | 0.49 | 0.56 | 0.44 | 0.51 | 0.49 | 0.93 | C |
| -0.9 | 29 | 0.03 | 0.10 | 0.97 | 0.03 | 0.90 | 0.10 | 0.13 | C |
| 0.9  | 30 | 0.00 | 0.00 | 0.00 | 1.00 | 0.00 | 1.00 | 0.00 | D |
| 0.4  | 30 | 0.19 | 0.20 | 0.19 | 0.81 | 0.20 | 0.80 | 0.39 | D |
| 0.0  | 30 | 0.36 | 0.36 | 0.36 | 0.64 | 0.36 | 0.64 | 0.72 | D |
| -0.4 | 30 | 0.46 | 0.48 | 0.54 | 0.46 | 0.52 | 0.48 | 0.94 | C |
| -0.9 | 30 | 0.05 | 0.06 | 0.95 | 0.05 | 0.94 | 0.06 | 0.11 | C |
| 0.9  | 31 | 0.00 | 0.00 | 0.00 | 1.00 | 0.00 | 1.00 | 0.00 | D |
| 0.4  | 31 | 0.20 | 0.19 | 0.20 | 0.80 | 0.19 | 0.81 | 0.39 | D |
| 0.0  | 31 | 0.36 | 0.36 | 0.36 | 0.64 | 0.36 | 0.64 | 0.72 | D |
| -0.4 | 31 | 0.47 | 0.45 | 0.53 | 0.47 | 0.55 | 0.45 | 0.92 | C |
| -0.9 | 31 | 0.07 | 0.05 | 0.93 | 0.07 | 0.95 | 0.05 | 0.12 | C |
| 0.9  | 32 | 0.00 | 0.00 | 0.00 | 1.00 | 0.00 | 1.00 | 0.00 | D |
| 0.4  | 32 | 0.21 | 0.17 | 0.21 | 0.79 | 0.17 | 0.83 | 0.38 | D |
| 0.0  | 32 | 0.36 | 0.36 | 0.36 | 0.64 | 0.36 | 0.64 | 0.72 | D |
| -0.4 | 32 | 0.49 | 0.44 | 0.51 | 0.49 | 0.56 | 0.44 | 0.93 | C |
| -0.9 | 32 | 0.09 | 0.03 | 0.91 | 0.09 | 0.97 | 0.03 | 0.12 | C |
| 0.9  | 33 | 0.01 | 0.00 | 0.01 | 0.99 | 0.00 | 1.00 | 0.01 | D |
| 0.4  | 33 | 0.22 | 0.16 | 0.22 | 0.78 | 0.16 | 0.84 | 0.38 | D |
| 0.0  | 33 | 0.36 | 0.36 | 0.36 | 0.64 | 0.36 | 0.64 | 0.72 | D |
| -0.4 | 33 | 0.50 | 0.41 | 0.50 | 0.50 | 0.59 | 0.41 | 0.91 | C |
| -0.9 | 33 | 0.14 | 0.02 | 0.86 | 0.14 | 0.98 | 0.02 | 0.16 | C |
| 0.9  | 34 | 0.01 | 0.00 | 0.01 | 0.99 | 0.00 | 1.00 | 0.01 | D |
| 0.4  | 34 | 0.23 | 0.15 | 0.23 | 0.77 | 0.15 | 0.85 | 0.38 | D |
| 0.0  | 34 | 0.36 | 0.36 | 0.36 | 0.64 | 0.36 | 0.64 | 0.72 | D |
| -0.4 | 34 | 0.51 | 0.39 | 0.49 | 0.51 | 0.61 | 0.39 | 0.90 | C |
| -0.9 | 34 | 0.18 | 0.01 | 0.82 | 0.18 | 0.99 | 0.01 | 0.19 | C |
| 0.9  | 35 | 0.01 | 0.00 | 0.01 | 0.99 | 0.00 | 1.00 | 0.01 | D |
| 0.4  | 35 | 0.23 | 0.13 | 0.23 | 0.77 | 0.13 | 0.87 | 0.36 | D |
| 0.0  | 35 | 0.36 | 0.36 | 0.36 | 0.64 | 0.36 | 0.64 | 0.72 | D |
| -0.4 | 35 | 0.53 | 0.36 | 0.47 | 0.53 | 0.64 | 0.36 | 0.89 | C |
| -0.9 | 35 | 0.21 | 0.00 | 0.79 | 0.21 | 1.00 | 0.00 | 0.21 | C |
| 0.9  | 36 | 0.02 | 0.00 | 0.02 | 0.98 | 0.00 | 1.00 | 0.02 | D |
| 0.4  | 36 | 0.25 | 0.10 | 0.25 | 0.75 | 0.10 | 0.90 | 0.35 | D |
| 0.0  | 36 | 0.36 | 0.36 | 0.36 | 0.64 | 0.36 | 0.64 | 0.72 | D |
| -0.4 | 36 | 0.53 | 0.31 | 0.47 | 0.53 | 0.69 | 0.31 | 0.84 | C |
| -0.9 | 36 | 0.25 | 0.00 | 0.75 | 0.25 | 1.00 | 0.00 | 0.25 | C |
| 0.9  | 37 | 0.03 | 0.00 | 0.03 | 0.97 | 0.00 | 1.00 | 0.03 | D |
| 0.4  | 37 | 0.25 | 0.06 | 0.25 | 0.75 | 0.06 | 0.94 | 0.31 | D |
| 0.0  | 37 | 0.36 | 0.35 | 0.36 | 0.64 | 0.35 | 0.65 | 0.71 | D |
| -0.4 | 37 | 0.55 | 0.20 | 0.45 | 0.55 | 0.80 | 0.20 | 0.75 | C |
| -0.9 | 37 | 0.29 | 0.00 | 0.71 | 0.29 | 1.00 | 0.00 | 0.29 | C |

Table 14: RP: realized payoff of Player 1; T1: Type 1 Error; T2: Type 2: Error; Error: Inferential Error; P2: Player 2

## Appendix C Opponent's Distribution

Figures 6 and 7 give the distribution of payoffs for player 2 for different values of realized payoff for player 1 under different correlation structures. For each scenario, two distributions are generated assuming that player 1 is cooperating: (1) the distribution of player 2 realized payoffs when she cooperates as well ( $CC$ ), and (2) the distribution of player 2 realized payoffs when she defects ( $CD$ ). One million simulations are done for each. For each simulation, all possible realized payoff for player 2 for a specific realized payoff for player 1 is plotted. A simulation is done for each distribution under every scenario. While we only display realized payoff of 29, 30, 31, 32, 33 (Figure 6) (we used the corresponding region for player 1 defecting: 41, 42, 43, 44, 45 (Figure 7) for player 1, we simulated all possible payoffs in the region of uncertainty including the boundaries.

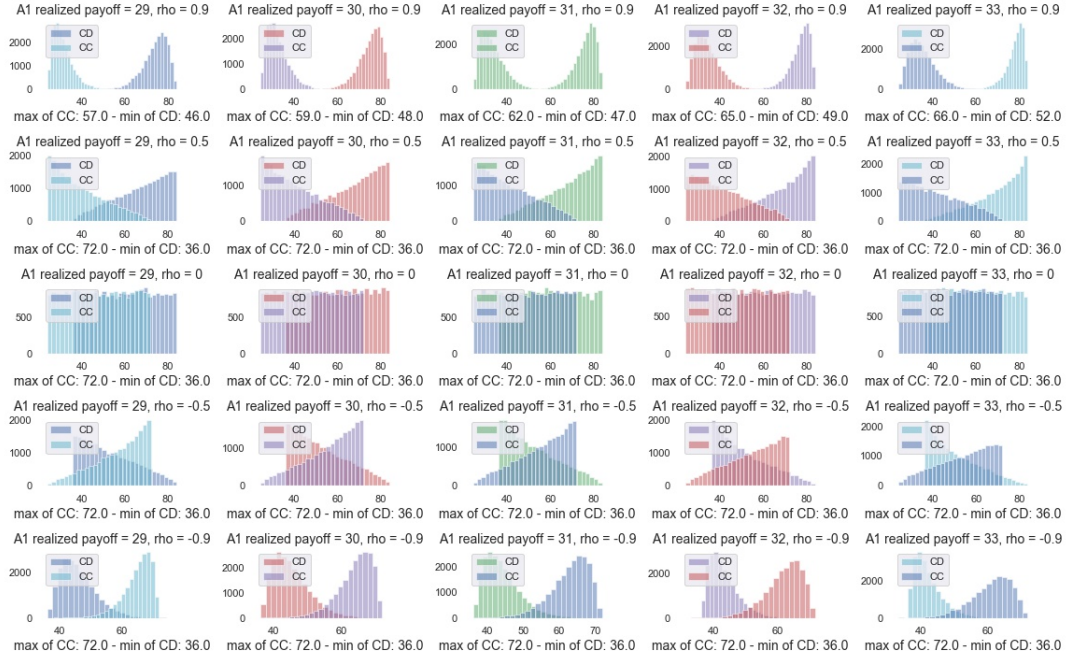


Figure 6: The distribution of player 2's realized payoff for various values of  $\rho$  and realized payoff of player 1, when player 1 plays cooperate.



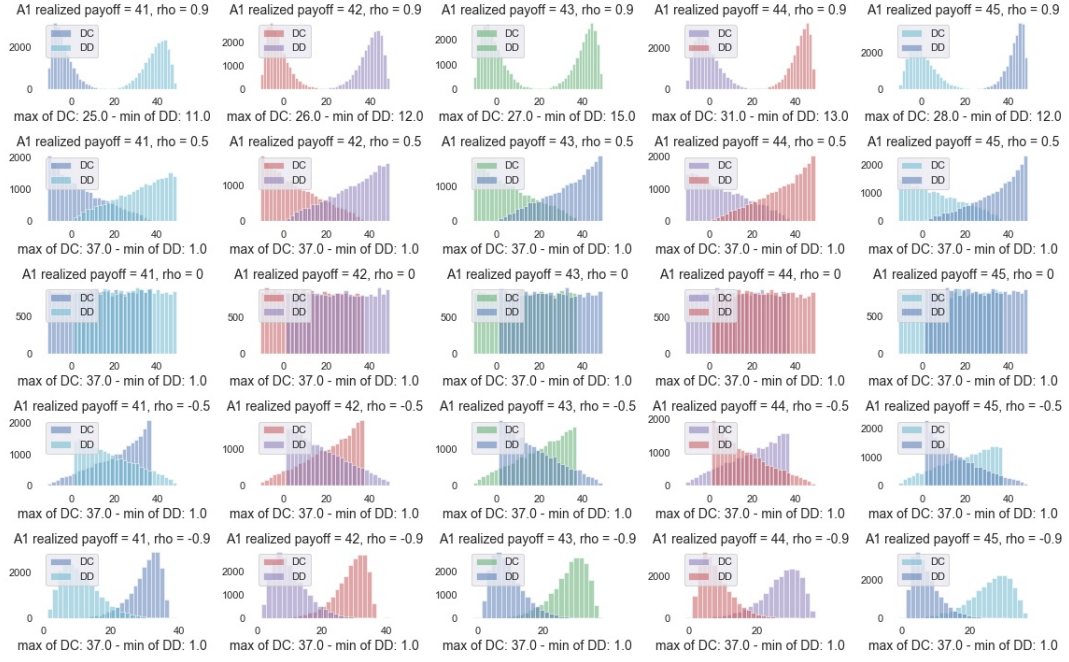


Figure 7: The distribution of player 2's realized payoff for various values of  $\rho$  and realized payoff of player 1, when player 1 plays defect.

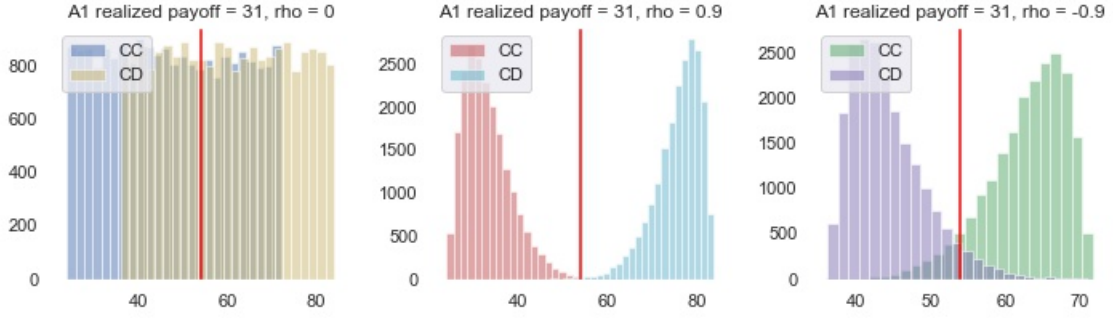


Figure 8: The distribution of player 2's realized payoff for a realized payoff of 31 for player 1, along with the benchmark value set by player 1. Player 1 will always get a signal that tells if player 2's realized payoff lies above or below this benchmark value. Player 1 is cooperating.

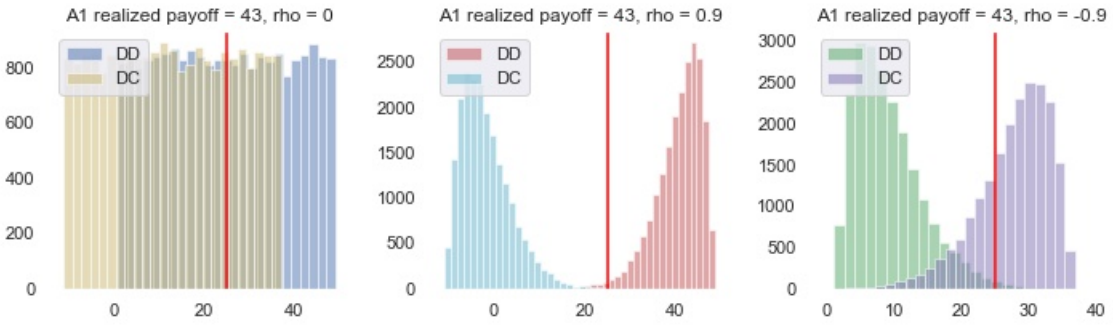


Figure 9: The distribution of player 2's realized payoff for a realized payoff of 43 for player 1, along with the benchmark value set by player 1. Player 1 will always get a signal that tells if player 2's realized payoff lies above or below this benchmark value. Player 1 is defecting.

Figure 8 gives an example of the signal player 1 receives about the realized payoff for player 2. Player 1 sets a benchmark value (the red vertical line) such that Type 1 and Type 2 errors are the same when  $\rho = 0$ . This benchmark value will also be used for every other values of  $\rho$ . Player 1 will receive a signal that tells if player 2's realized payoff is above or below this benchmark value. If we focus on the scenario where player 1 has a realized payoff of 31, if shocks are positively highly correlated, the mass of the  $CC$  distribution lies below the benchmark value, and the mass of the of the  $CD$  distribution lies above the benchmark value. If player 1 is told that player 2's realized payoff is above the benchmark value (in other words, player 2 got a relatively big realized payoff), if they assumed player 2 defected, the probability of being correct is very large (area of the  $CD$  distribution above the benchmark value) and the probability of being wrong (area of the  $CC$  distribution above the benchmark value) is very small. The probability of being incorrect increases as  $\rho$  decreases. The analysis is similar if player 1 is defecting. Figure 9 shows the distributions when player 1 defects.

## Appendix D Experimental Instructions

### Instructions for the Noise Treatment

#### Welcome!

Today's experiment will last about 60 minutes. You will be paid a show-up fee of \$5 together with any money you accumulate during this experiment. The money you accumulate will **depend partly on your actions, partly on the actions of others, and partly on chance. Therefore, please read the instructions carefully.** This money will be paid at the end of the experiment in private and in cash.

Your returns will be recorded in points. At the end of the session, the total number of points in your account will be converted into cash at an exchange rate of 300 unit = \$1. It is possible for you to get negative points in a round. If at the end of the session you have negative units in your account, you will be paid the show-up fee.

It is important that during the experiment you remain silent. If you have any questions or need assistance of any kind please raise your hand, but do not speak, and an experiment administrator will come to you and you may then whisper your question.

In addition, please turn off your cell phones and put them away now. Please do not look into anyone's booth at any time.

Please read the following instructions carefully. You will be given a quiz at the end to test your understanding and you earn \$0.50 for each correct answer.

#### Agenda:

- Experimental instructions
- Quiz
- Experiment

#### How a match works

This session is made up of 10 matches between you and other participants in the room. In each match, you play a random number of rounds with another participant. The length of a match is randomly determined in that, after each round, there is a 90% chance that the match will continue for at least another round. Once the match ends, you will be randomly re-grouped with another participant to play another match. Whenever a match ends, you will be informed of this before you are re-grouped.

#### Decisions and Payoffs (Before Random Draw)

In each match, you will make a series of investments in a project with the same participant. For each round of a match, you can invest in either Project A or Project B. The participant you are playing with has the same two options. You each choose your project at the same time. The returns

on investment depends on the project you choose, the project the other participant chooses and a random draw. That is, the returns you get depend on:

- the investment you made (Project A or Project B)
- the investment made by the other participant
- a random draw

The following table summarizes the return you get based on your decision and the other participant's decision:

|             |   | Other participant's choice |              |
|-------------|---|----------------------------|--------------|
|             |   | A                          | B            |
| Your choice | A | <b>48,48</b>               | <b>13,60</b> |
|             | B | <b>60,13</b>               | <b>25,25</b> |

The first red bolded entry in each cell represents your returns before accounting for the random draw, while the second entry in blue represents the returns of the participant you are grouped with (how the random draw affects you and the other participant's returns will be explained below). Ignoring the random draw, if:

- You invest in Project A and the other participant invests in Project A, you both earn **48 points**
- You invest in Project A and the other participant invests in Project B, you earn **13 points** and the other participant earns **60 points**
- You invest in Project B and the other participant invests in Project B, you both earn **25 points**
- You invest in Project B and the other participant invests in Project A, you earn **60 points** and the other participant earns **13 points**

### **Your project returns may change depending on a random draw**

In each round, after you have invested in a project, your return may change by a random draw. Let's call this random draw  $v_1$ . This means that, your return may increase, decrease, or stay the same by an amount  $v_1$ . The computer will randomly select this number in each round. This random draw does not depend on the project that you or the other participant choose or the random draw in previous rounds. This draw is completely random.

This random draw will always be a number from -24 to 24. Each number, of the 49 integer values between -24 and 24, is equally likely to occur.

The other participant's return, after they have invested, will also change by a random draw. Let's call this amount  $v_2$ . This means that, the returns from their project may increase, decrease, or stay the same by an amount  $v_2$ . The computer will generate these integers for you both. **These integers will**

be completely independent. This means that your random draw is completely unrelated to the random draw of the other participant.

In the diagram below, there are 500 examples of random draws for you and the other participant, where each dot represents a random draw. If you hover your cursor over one of these dots, you will see a pair of numbers where the first integer (labelled 'yours') represents your random draw and the second integer (labelled 'theirs') is the random draw for the other participant.



Now, the total return you receive is dependent on your random draw AND the choices made by the other participant. The other participant's random draw does not affect your return. With the random draw, the rule for your investment returns now becomes:

|             |          | Other participant's choice |                      |
|-------------|----------|----------------------------|----------------------|
|             |          | <i>A</i>                   | <i>B</i>             |
| Your choice | <i>A</i> | $48 + v_1, 48 + v_2$       | $13 + v_1, 60 + v_2$ |
|             | <i>B</i> | $60 + v_1, 13 + v_2$       | $25 + v_1, 25 + v_2$ |

Pay close attention to the following information.

For both you and the other participant, taking into account the return and the random draw:

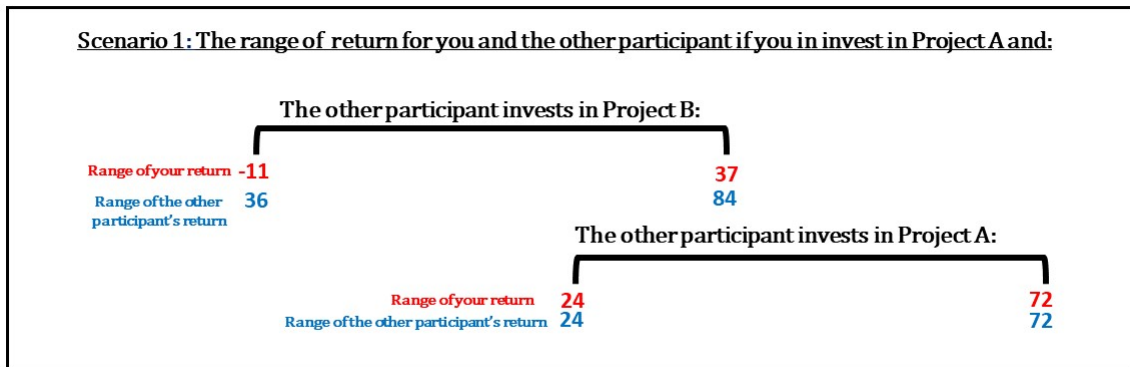
- the minimum return that can be received is **-11** (the lowest possible return 13 plus minimum possible random draw of -24)
- and the maximum return that can be received is **84** (the highest possible return 60 plus maximum possible random draw of +24).

After you have made an investment choice in A or B, the range of the returns after accounting for the random draw is:

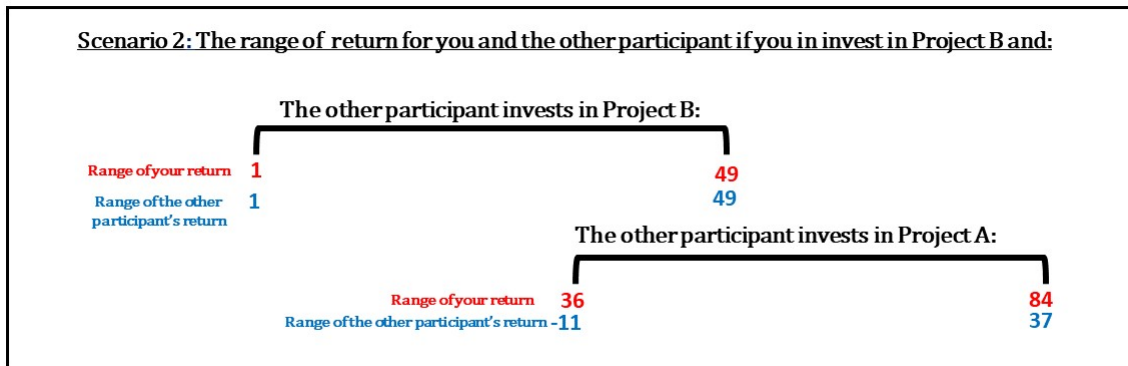
1. If you invest in **Project A** and the other participant invests in **Project A**:

- Your return will range from **24** (48 plus worst random draw -24) to **72** (48 plus best random draw +24)

- The other participant's return will range from **24** (48 plus worst random draw -24) to **72** (48 plus best random draw +24)
2. If you invest in Project A and the other participant invests in Project B:
    - Your return will range from **-11** (13 plus worst random draw -24) to **37** (13 plus best random draw 24)
    - The other participant's return will range from **36** (60 plus worst random draw -24) to **84** (60 plus best random draw +24)
  3. Note that if YOUR return falls between 24 and 37, you CANNOT know for sure whether the other participant invested in Project A or B. If YOUR return is less than 24 or greater than 37 you CAN know for sure what project the other participant invested in. (see Scenario 1 in the diagram below).



1. If you invest in Project B and the other participant invests in Project B:
  - Your return will range 1 to 49
  - The other participant's return will range from 1 to 49
2. If you invest in Project B and the other participant invests in Project A:
  - Your return will range 36 to 84
  - The other participant's return will range from -11 to 37
3. Note that if YOUR return falls between 36 and 49, you CANNOT know for sure whether the other participant invested in Project A or B. If YOUR return is less than 36 or greater than 49 you CAN know for sure what project the other participant invested in. (see Scenario 2 in the diagram below).



### You will get a signal

You will not be told the return of the other participant, but you will always get a signal about their return. This signal will tell if the other participant's return is above, below, or equal to a benchmark value. This benchmark value can help you rule out ranges of values of the other participant's return. You will set one benchmark value that will be used if you select Project A and another if you select Project B. These benchmark values have to be within the range of possible return for both projects for the other participant. That is, 24 to 84 for Project A and -11 to 49 for Project B.

At the beginning of each match, you will be prompted to select these benchmark values. For example, if you set a benchmark value of 50 for Project A, you will be signaled that the other participant's return is above, below, and equal to 50. You only set this benchmark value at the beginning of each match. The same benchmark will be used for all the rounds in a match. When a new match begins, you will be prompted to set the benchmark value again.

Here you can simulate how the benchmark can be used. First select a project, then the computer will randomly select a project as well. Move the slider to see what information you receive about the other participant based on the benchmark you select.

Practice for using the Benchmark [Please take a few minutes to try this!]

Click on a project then use the slider below to see how the information you get about the other participant changes with all possible benchmark values. You can click as many times as you want.

In each cell, the amount to the left and **bolded in red** is the return for you, and the one to the right in **blue** is for the other participant.

|     |           | The Other Participant |                |
|-----|-----------|-----------------------|----------------|
|     |           | A                     | B              |
| You | Project A | <b>48</b> , 48        | 13 , <b>60</b> |
|     | Project B | <b>60</b> , 13        | <b>25</b> , 25 |

**You selected Project A and as a result your return is 23**

This slider represents the different benchmark values you can choose

Click to select a benchmark value. Remember to select a project first!!



If you select a benchmark value of 76, you will be told that: The other participant's return is less than your benchmark value

### After Each Round

On the result page, **we will ask you what you think the chances are that the other participant chose Project A.**

Depending on your guess, you can earn 2 points or 0 points. We are interested in learning about your best and honest guesses. **You will be paid according to a formula which is specifically designed to maximize the chances that you will win the 2 points if you submit your best guess.**

**Your guess will be converted into a chance-to-win.** This is calculated by the computer according to a formula that is explained on separate page that you can request after the experiment. On the computer interface, you will be able to see the chance-to-win for each outcome directly below your guess.

You will not be paid for your answer until the end of the experiment. Your answer will not be shown to any other participant. Your answer will not affect the experiment in any way.

### The Interface of the Experiment

Before each match, your computer screen will look like this:



**Please select your benchmark values.**

You will select two values:

- The benchmark value you want to use if you select Project A (an integer between 24 and 84)
- The benchmark value you want to use if you select Project B (an integer between -11 and 49)

What is your benchmark value when you select Project A:

What is your benchmark value when you select Project B:

Figure 10: Set Your Benchmark Values

After you select these benchmark values, the round will begin. In each round, the screen to select a project for you and the other participant looks like this:

The benchmark values you selected at the beginning of the match

## Choose Your Project

A summary of outcome for previous rounds in Match 1:  
 My Benchmark value if I choose Project A: 64  
 My Benchmark value if I choose Project B: 10

Your history up to this round

| Round | Project | My Return | Other's return to Benchmark |
|-------|---------|-----------|-----------------------------|
| 1     | B       | 55 points | Above                       |
| 2     | B       | 11 points | Above                       |
| 3     | B       | 39 points | Above                       |

Select Project A or Project B

Your return may change by a random draw after you have invested in a project. Whatever random draw you face is completely independent of the random draw faced by the other participant.

In each cell, the amount to the left and **bolded in red** is the return for you, and the one to the right in **blue** is for the other participant.

|     |           | The Other Participant |         |
|-----|-----------|-----------------------|---------|
|     |           | A                     | B       |
| You | Project A | <b>48</b> , 48        | 13 , 60 |
|     | Project B | 60 , 13               | 25 , 25 |

Figure 11: SYour Selection Screen

This screen also displays a summary of the outcome of all previous rounds including: the project you chose; your return (inclusive of the random draw you faced); and if the other participant's return is above, below or equal to the benchmark value you set at the beginning of the match.

After you and the other participant have made a decision, your result screen will display:

- The decision made BY YOU
- YOUR realized returns (inclusive of your random draw).
- If the other participant is above, below or equal to the benchmark
- A slider for you to guess the probability that the other participant selected Project A

This is an example of what the computer screen may look like after you have made your choice:

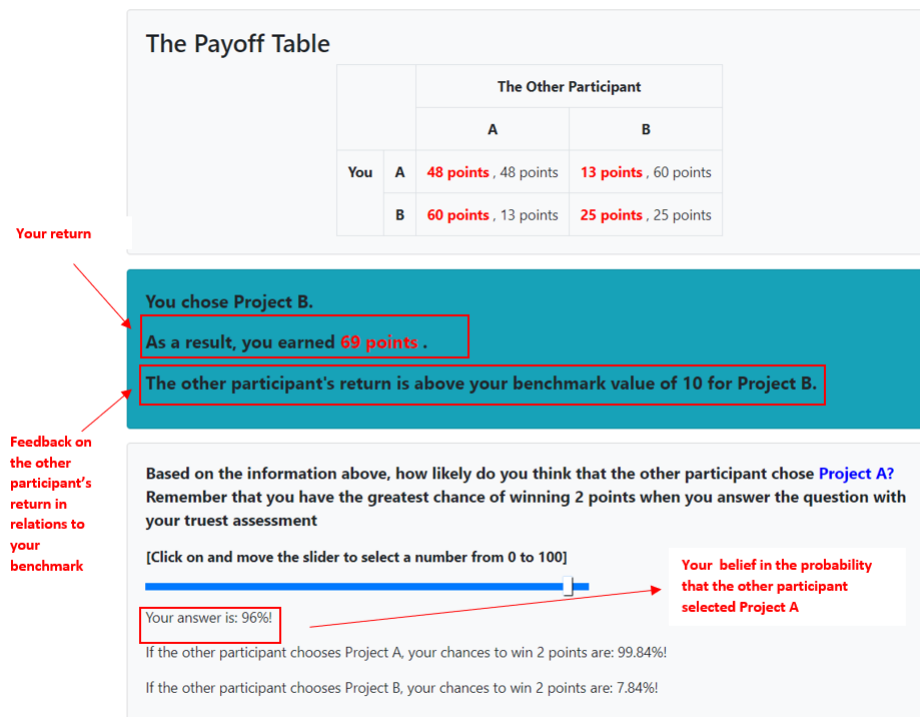


Figure 12: Your Result Screen

Once a match ends, you will be randomly re-grouped with a different participant in the room for another match. Each match has the same setup. You will play a number of such matches with different people in the room.

### **Reminders**

To summarize, the number of rounds in a match is randomly determined. After each round, there is a 90% chance that the match will continue for at least another round. You and the other participant will get a random draw. Whatever random draw you get is completely independent of the random draw faced by the other participant. This means that your random draw is completely unrelated to the random draw of the other participant.

You will not know the return of the other participant, but you will be able to select a benchmark value to signal to you the possible range of returns for the other participant. After you both have invested in a project, you will be told if the other participant's return was above, below or equal to this benchmark value. At the end of this session, you will receive \$1 for every 300 point in your account. You will now take a very short quiz to make sure you understand the setup. You will earn \$0.50 for each correct answer.

After the quiz, you will play 4 practice rounds to get you familiarized with the game. For the practice rounds, you will play against the computer and NOT the other participants. The computer will randomly select responses. Also, you will be able to select benchmark values for each round. This **ONLY** happens for the practice rounds. After the practice rounds, you will begin playing with the other participants in the room.

## Appendix E Description of the Belief Elicitation strategy

For this, you will receive either 0 points or 2 points. Your chance to win 2 points depends on both your guess and if the other participant invested in Project A. Specifically, your chance of receiving 2 points is determined in the following way:

1. First, you will guess the probability that the other participant invested in Project A. You will guess a number from 0 to 100, that we convert to a decimal.
2. If the other participant invested in Project A, your chance-to-win 2 points will be:

$$2z - z^2$$

where  $z$  is the probability you selected, that the other participant selected Project A

3. If the other participant invested in Project B, your chance-to-win 2 points will be:

$$1 - z^2$$

4. To determine whether you receive 2 points, the computer will randomly draw a number between 0 and 100. Each number between 0 and 100 is equally likely to be picked
5. If the number drawn by the computer is less than or equal to your chance-to-win, then you will receive 2 points. Otherwise, you receive 0 points

## Appendix F Feedback After Supergame

Match 1 has ended. Your cumulative decision and payoff for this match is:

All periods:

| Rounds | Your Decision | Your Return | Other Decision | Other Return |
|--------|---------------|-------------|----------------|--------------|
| 1      | A             | -8 points   | B              | 39 points    |
| 2      | B             | 38 points   | A              | -10 points   |
| 3      | A             | 35 points   | B              | 82 points    |
| 4      | B             | 66 points   | A              | 20 points    |
| 5      | B             | 51 points   | A              | -5 points    |
| 6      | A             | 50 points   | A              | 52 points    |
| 7      | A             | 41 points   | A              | 39 points    |
| 8      | B             | 53 points   | A              | 17 points    |
| 9      | A             | 65 points   | A              | 60 points    |

## Appendix G The Maximum Likelihood Estimation Method

We use the Strategy Frequency Estimation Method (SFEM) from Dal Bó and Fréchette (2011) to estimate the fraction of strategies employed in each treatment. This methodology uses a Maximum Likelihood Estimation (MLE) to estimate the frequency with which each strategy from a set of pre-determined set of strategies is found experimental data. This methodology has since been employed Fudenberg, Rand, and Dreber (2012), Rand, Fudenberg, and Dreber (2015), Dal Bó and Fréchette (2018), Aoyagi, Bhaskar, and Fréchette (2019), Dal Bó and Fréchette (2019) and Romero and Rosokha (2019), for example. This method assumes that each subject uses the same strategy across supergames. However, they can make mistakes. These mistakes are not the errors that are generated from the experimental design, but rather, it is assumed that subjects can make mistakes when choosing their intended actions for the particular strategy they are following.

Using the notations of Dal Bó and Fréchette (2011), assume that the probability with which subject  $i$  makes mistakes is  $1 - \beta$  and the probability that her chosen actions correspond with a strategy  $k$  is  $\beta$ . The likelihood that her observed choices were actually generated by strategy  $k$  is  $Pr_i(s^k) = \prod_{M_i} \prod_{R_{im}} (\beta)^{I_{imr}^k} (1 - \beta)^{1 - I_{imr}^k}$ .  $I_{imr}^k$  is an indicator function that takes the value 1 when the choice that was actually made in round  $r$  and supergame  $m$  is the same as what the subject would have made if she were following strategy  $k$ . It is coded 0 otherwise.  $M$  and  $R$  are the sets of supergames

and rounds.  $\beta$  is estimated within the model. It can also be interpreted as the probability that an action is taken given that it is prescribed by a strategy  $k$ . Therefore  $\beta$  is the basis for evaluation of model fit, that is, as the model fit improves  $\beta$  approaches 1.

Therefore, the MLE process entails choosing both the probability of mistakes and the frequency of strategies that maximizes the likelihood of the sequences of choices. That is, the log-likelihood is  $\sum_I \ln(\sum_K \phi^k Pr_i(s^k))$ , where  $K$  is the subset of strategies being considered and  $\phi^k$  is a vector of parameter estimates that represent the frequency of strategies.

We bootstrapped the standard errors in a way that respects the data generating process of our experimental data. We randomly draw the appropriate number of sessions, then for each session the appropriate number of subjects, then supergames. All with replacement. The bootstrapping process was done 1000 times. The standard deviation of the bootstrapped MLE estimates provide the standard errors.

## Appendix H Description of Strategies

| Strategy | Description   |
|----------|---|
| AD       | Always defect   |
| DTFT     | Defect in the first round, then play TFT  |
| DGRIM2   | Defect in the first round, then play GRIM2  |
| GRIM     | Cooperate until the other player defects, then defect forever                         |
| TFT      | Cooperate unless other player played defection in the last round                      |
| 2TFT     | Cooperate unless other player defected in either of the last two rounds               |
| GRIM2    | Cooperate until the other player defects in 2 consecutive rounds, then defect forever |
| GRIM3    | Cooperate until the other player defects in 3 consecutive rounds, then defect forever |
| TF2T     | Cooperate unless other player defected in both of the last two rounds                 |
| AC       | Always cooperate  |

## Appendix I Details on the Automata

There are 128 combinations of automata using 7-bit strings ( $2^7$ ). However, different automata can represent the same strategy. As such, there are 26 unique strategies. Of the 128 strategies, forty of them represents AD and forty represents AC. The other 24 unique strategies have two different 7-bit string representation.

To facilitate the noisy signal – two benchmark values – like subjects in the lab, strategies select these at the beginning of each generation. While we cannot theorize the process through which subjects select these benchmark values, the automata are programmed to randomly select these. With this method, we can see what benchmark values survive the evolutionary process. To accommodate the strategy and benchmark values, a 21-bit string is generated for each strategy. In Figure 13, we show an example of a representation for TFT. The first 7-bit string translate into a benchmark value of 44 if C is played. The second 7-bit string translated into 15 if D is played. The third 7-bit string represents the strategy.

Figure 13: TFT Representation

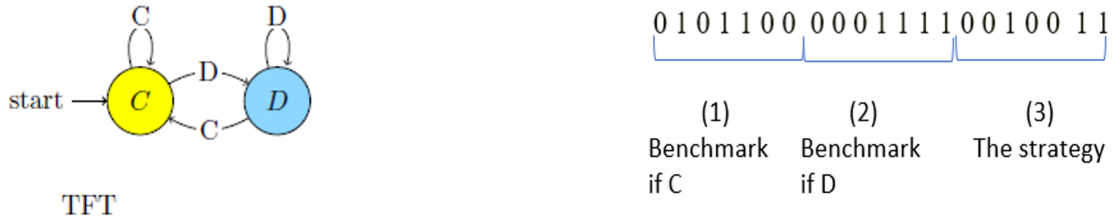


Figure 14 shows how the strategies are coded along with the representation for TFT and DTFT. The strategies are coded as having an initial state and internal states. The first bit indicates that the automaton begins in state 0. The other two bits (first bit for state 0 and second bit for state 1) prescribes the action for each state. The next two gives the transition rule if cooperation is observed in each state and the final two bits give the transition for each state if defection is observed. For example, for TFT, the final 6 bits suggest the following action. The first two [0 1] says play C in state 0 and play D in state 1. The second two [0 0] prescribes the action in each state if the other player is observed to have cooperated. It says that if you are in state 0, transition to state 0. And, if you are in state 1, transition to state 0. The final 2 bits [1 1] prescribes the action in each state if the other player is observed to have defected. If you are in state 0 transition to state 1 and if you are in state 1 transition to state 1. Note that, in the diagrammatic representation, the transition states are represented by the labelled arcs and a vortex shows the internal state. The player's prescribed action is represented by the letter in the middle.



Figure 14: Coding of the Strategies

| <b>X</b><br>Starting<br>State   | <b>XX</b><br>The action<br>prescribed for the  | <b>X X</b><br>In each state, what<br>state to transition<br>to if C is observed  | <b>X X</b><br>In each state,<br>what state to<br>transition to if D<br>is observed   |
|---------------------------------|--|--|--|
| <b>TFT</b>                      |  |  |  |
| <b>0</b><br>Start in<br>State 0 | <b>0 1</b><br>C in state 0<br>and D in state   | <b>0 0</b><br>In state 0 if C is observed,<br>transition to state 0. In state 1, if<br>C is observed, transition to state<br>0 | <b>1 1</b><br>In state 0 if D is observed,<br>transition to state 1. In state<br>1, if D is observed, transition<br>to state 1 |
| <b>DTFT</b>                     |  |  |  |
| <b>0</b><br>Start in<br>State 0 | <b>1 0</b><br>D in state 0 and C in<br>state 1 | <b>1 1</b><br>In state 0 if C is observed,<br>transition to state 1. In state 1, if C<br>is observed, transition to state 1    | <b>0 0</b><br>In state 0 if D is observed,<br>transition to state 0. In state 1, if D<br>is observed, transition to state 0    |

Figure 15: The 26 Unique Automata

