

# Revealed Reputations in the Finitely-Repeated Prisoners' Dilemma

Caleb A. Cox · Matthew T. Jones · Kevin E. Pflum ·  
Paul J. Healy

November 19, 2014

**Abstract** In a sequential-move, finitely-repeated prisoners' dilemma game (FRPD), cooperation can be sustained if the first-mover believes her opponent might be a behavioral type who plays a tit-for-tat strategy in every period. We test this theory by revealing second-mover histories from an earlier FRPD experiment to their current opponent. Despite eliminating the possibility of reputation building, aggregate cooperation actually increases when histories are revealed. Cooperative histories lead to increased trust, but negative histories do not cause decreased trust. We develop a behavioral model to explain these findings.

**Keywords** prisoners' dilemma · finitely-repeated games · cooperation · reputation-building  
**JEL** C70 C73 C92

## 1 INTRODUCTION

Cooperation in the finitely-repeated prisoner's dilemma (FRPD) is both widely observed and difficult to rationalize. The model developed by Kreps, Milgrom, Roberts, and Wilson (1982)—the most prominent theory to justify this behavior—shows that such cooperation can be rational if one player believes her opponent might be a “behavioral type” who plays the tit-for-tat strategy

---

Caleb A. Cox  
Durham University Business School  
Mill Hill Lane  
Durham DH1 3LB, U.K  
E-mail: caleb.cox@durham.ac.uk

Matthew T. Jones  
Federal Trade Commission  
600 Pennsylvania Avenue NW, Mail Drop HQ-238  
Washington, D.C. 20580, U.S.A.  
E-mail: mjone1S@ftc.gov

Kevin E. Pflum, *corresponding author*  
University of Alabama  
Box 870224  
Tuscaloosa, AL, 35487, U.S.A.  
205-348-0788  
E-mail: kpflum@cba.ua.edu

Paul J. Healy  
The Ohio State University  
1945 North High Street  
Columbus, Ohio 43210, U.S.A.  
E-mail: healy.52@osu.edu

regardless of the history of play. The opponent can then take advantage of these beliefs by imitating the behavioral type early in the game and defecting as the final period approaches, resulting in a pattern of cooperation consistent with aggregate-level data from laboratory experiments (e.g., Andreoni and Miller, 1993).

This reputation-building theory requires that players have some uncertainty about their opponents' types. In our experiment, however, we find that cooperative play persists even when reputation building is rendered impossible by revealing players' histories of play. Specifically, we have subjects who completed a block of five sequential-move FRPD games against varying opponents, play a second block of five games against new opponents, who can see the subjects' histories of play from the first block.<sup>1</sup> Selfish, rational players cannot credibly imitate the behavioral type in the second block because their true colors have been revealed through their history of play in the first block. Despite eliminating type uncertainty, we find aggregate patterns of cooperation very similar to treatments where no history is revealed. On the individual level, first-movers who are relatively distrusting (seldom cooperating in the first block) tend to be more cooperative in the second block, even when the second-mover's revealed history is relatively uncooperative. Hence, rather than reducing cooperation by eliminating the opportunity for reputation building, revealing histories of play generally improves cooperation. This finding is clearly inconsistent with standard reputation-building explanations of cooperation in FRPDs.

We organize these results through a model of semi-rational behavior similar to those of Kreps et al (1982) and Radner (1986). Here, players decide in which round to stop playing tit-for-tat and begin unconditionally defecting by weighing the risk of cooperation against the immediate gain from defecting. Unlike Kreps et al (1982), players in our model form arbitrary or "naïve" prior beliefs about how many rounds their opponents will continue playing tit-for-tat, which may be inconsistent with the opponent's actual strategy. Players decide how long to conditionally cooperate in each game based only on these naïve prior beliefs and information about their opponents' history of play, if revealed. Because this model does not assume any higher-level reflection about the rationality or best-response of the opponent, it provides a contrasting benchmark to a model of full, commonly-known rationality.

One implication of this model is that cooperation is sustainable for many rounds even when players have relatively pessimistic beliefs about their opponents' strategies. For example, the model predicts that cooperation will be sustained until the penultimate round if the players have a uniform prior; i.e., if players believe their opponent will defect with equal probability in every round of the FRPD. In stark contrast to the reputation-building theory, the model also predicts that the level of cooperation will be the same or even higher when players learn that their opponent is not a behavioral type. This pattern of behavior is frequently observed in our experiment and sufficient to reject reputation-building as an explanation for cooperation. Hence, we find that this semi-rational model of naïve beliefs rationalizes the observed experimental behavior.

The paper is organized as follows. In section II, we review the related theoretical and experimental literature. Section III contains a description of the reputation building theory of cooperation in FRPDs. Section IV details the experimental design and the results of the experiment are presented in Section V. In Section VI, we propose a model of boundedly-rational cooperation that

---

<sup>1</sup>During the first block, subjects know that there will be an optional second experiment but know nothing about its nature.

explains the experimental results. Section VII concludes with a summary of the main results. The Appendix contains the proofs and additional summary statistics.

## 2 RELATED LITERATURE

Many experiments have studied the consistency of players' behavior in FRPDs with the theory of Kreps, Milgrom, Roberts, and Wilson (1982). For example, Andreoni and Miller (1993) compare the amount of cooperation in 10-round FRPDs with one-shot prisoners' dilemma games. They find significantly higher cooperation in the FRPDs compared to the one-shot games, as well as significantly higher cooperation in early rounds compared to later rounds. These patterns are consistent with the reputation-building theory of Kreps et al (1982) at an aggregate level, though in a similar FRPD experiment Cooper, DeJong, Forsythe, and Ross (1996) observe that, at the individual level, only 25% of subjects play consistently with reputation building. Cooper et al argue that the time path of play exhibits more cooperation than the Kreps et al model predicts and speculate that their findings could indicate reputation building if they were to consider alternative types of "irrational" players.<sup>2</sup>

Camerer and Weigelt (1988) study the reputation-building sequential equilibrium in a finitely-repeated investment game that is similar in structure to our FRPD game. The authors randomly assign a small fraction of second-movers to have payoffs such that they prefer cooperation over defecting. This exogenously induces the behavioral type. They find evidence consistent with the equilibrium prediction: as time progresses, first-movers are less likely to trust and second-movers are more likely to defect. The observed mixing probabilities are different than predicted by the induced probability of the behavioral type, but can be explained easily by assuming first-movers believe an additional 17% of second-movers with the "non-cooperative" incentives still prefer cooperation. To test this theory, the authors run additional experiments in which all second-movers were given non-cooperative payments. The sequential equilibrium prediction assuming beliefs of 17% is calculated, and the data from the new experiments conform surprisingly well to that prediction. Thus, Camerer and Weigelt (1988) provide strong evidence in favor of a reputation-building theory in which the first-mover's beliefs are "homemade" naturally, and need not be induced in the lab.

Neral and Ochs (1992) extend Camerer and Weigelt (1988) by analyzing the behavioral responses of players to changes in the parameters of the game. Like Camerer and Weigelt they find that uncertainty induces players to develop a mutually profitable relationship consistent with the predictions of sequential equilibrium. However, Neral and Ochs find problems with the comparative static predictions of the sequential equilibrium model. When the parameters of the game are altered (e.g., decrease the payoff to second-mover) they find that the players respond in the exact opposite direction from what the theory predicts and that these results cannot be explained using the homemade priors specified by Camerer and Weigelt.<sup>3</sup>

---

<sup>2</sup>In contrast with the finitely repeated case, experimental evidence has shown that cooperation in the infinitely repeated prisoners' dilemma aligns well with theoretical predictions. For example, Roth and Murnighan (1978) and Murnighan and Roth (1983) study behavior in indefinitely-repeated prisoners' dilemma experiments and find behavioral differences predicted by standard folk theorem equilibria. More recently, Dal Bó (2005) finds experimental evidence that greater cooperation occurs in an indefinitely repeated prisoner's dilemma with the same expected length as a finitely repeated control and Dal Bó and Frechette (2011) find evidence that subgame perfection is a necessary (but not sufficient) condition in supporting cooperation in an indefinitely repeated prisoners' dilemma.

<sup>3</sup>Jung, Kagel, and Levin (1994) analyze the sequential equilibrium of a chain-store game that shared some features with Camerer and Weigelt's borrower-lender game and also find discrepancies with the theory that cannot be resolved with

More recently, Reuben and Suetens (2012) find that many players cooperate in a manner consistent with Kreps et al (1982). Reuben and Suetens use the strategy method to disentangle strategically and non-strategically motivated cooperation in a sequential prisoner's dilemma with an uncertain endpoint in which cooperation is not an equilibrium strategy for rational players.<sup>4</sup> Players in their experiment can condition their action on whether they are currently playing the last period of the game or whether the game will continue. Second-movers who cooperate as long as the first-mover cooperates unless it is the final round are classified as reputation builders. Among second-movers who cooperate, the authors find that a third to two-thirds do so for reputation-building rather than reciprocity. Moreover, Reuben and Suetens find that an increase in the payoff of mutual cooperation increases the ratio of reputation builders to unconditional defectors, consistent with Kreps et al (1982).

Kagel and McGee (2014) compare individual play and team play in the finitely-repeated prisoners' dilemma. In the team play treatment, each player in the game is a 2-person team that can chat internally (but cannot communicate with opponents) and makes joint decisions. They find that teams are initially less cooperative than individuals, but with experience become more cooperative. Kagel and McGee's analysis of team chat logs suggests that cooperation is driven by a failure of common knowledge of rationality, as teams attempt to anticipate when their opponents might defect and try to defect one period earlier, without accounting for the possibility of their opponents thinking similarly. Inconsistent with Kreps et al (1982), they find no evidence that players anticipate opponents' beliefs and attempt to mimic an irrational player, while they find that increased chat about the round in which the opponent will defect is accompanied by an increase in cooperation over the course of the experiment. This finding is also consistent with Embrey et al (2014), who find that players learn to play "threshold strategies" in which they conditionally cooperate for a fixed length of time and then defect. Though players converge to these strategies, unraveling of cooperation occurs very slowly across supergames. Subjects clearly are not applying backwards induction reasoning.

Other experiments have examined how cooperation in one-shot prisoners' dilemma games is impacted when players see their opponents' play history. Schwartz, Young, and Zvinakis (2000), Camera and Casari (2009), and Gong and Yang (2010) all find that observing an opponent's history of play significantly increases cooperation, though Duffy and Ochs (2009) do not observe this effect in their data. In all of these studies, subjects are aware that their actions will be revealed to future opponents. Thus, players can follow a reputation-building strategy even though their opponents differ in every period. In contrast, because subjects in our experiment are not told in the first part of the experiment that their play histories will be revealed later, they are unlikely to perceive any benefit from cooperating in the final period of each repeated game in the first part.

Alternative models have been proposed that predict rational cooperation in FRPDs. For example, Selten and Stoecker (1986) develop an alternative theory based on a Markov learning model. In their model players establish a period of first defection and update their period of intended defection based on their experience in the previous supergame. The authors conduct an experiment in which subjects play 25 ten-period FRPD supergames and find that cooperation ends in

---

an appeal to homemade beliefs. Similarly, both Brandts and Figueras (2003) and Tingley and Walter (2011) find higher rates of cooperation than predicted by reputation building in shorter finitely repeated games.

<sup>4</sup>Reuben and Suetens ensure that cooperation is not a rational strategy by setting the probability that the game terminates below the threshold required for cooperation to constitute a subgame perfect equilibrium.

earlier periods as subjects gain experience.<sup>5</sup> Subjects in their experiment are told that they will play each opponent only once and are given no information about opponents' histories of play in prior supergames. Hence, their data do not address the validity of reputation building directly, nor are their results inconsistent with reputation building.<sup>6</sup> By revealing second-mover histories in one treatment and not in another, our design allows a more direct examination of how a player's intended period of first defection and beliefs about her opponent's intentions matter for cooperation.<sup>7</sup>

Other studies have focused on the role of reputation in inducing cooperation, though not in the context of an FRPD game. Gächter and Thoni (2005) and Ambrus and Pathak (2011) show how cooperation can be sustained in a public goods contribution game when some players are selfish and others are reciprocating with varying information on other players' past behavior. As in our design, Ambrus and Pathak incorporate a restart in their experimental design to see how it impacts cooperation, as do Gächter and Thoni (2005). In Gächter and Thoni (2005), the past behavior of subjects is revealed and the response to these "reputations" in the public goods game is studied. In Ambrus and Pathak, players in their experiment know in advance that they will be participating in subsequent games, but in Gächter and Thoni (2005), like our experiment, they do not. Bolton, Katok, and Ockenfels (2005) examine how information about their partner in an image scoring game affects cooperation, while Irlenbusch and Sliwka (2005) study the role of reputation and uncertainty about the partner's type in inducing cooperation in a gift-exchange setting. The former find that providing players with more information about their partner's last action as well as the action of their partner's previous partner increases cooperation while the latter find that direct reciprocal behavior is stronger when efforts are revealed. Healy (2007) considers a reputation-building equilibrium when firms stereotype workers and find that selfish workers imitate fair-minded types when firms have sufficiently high priors to generate cooperation. Similarly, Roe and Wu (2009) find evidence for the reputation-building equilibrium by finding that employees classified as selfish mimic cooperative employees when individual histories are observable, but not when histories are kept private. Andreoni and Croson (1998) survey repeated public goods game experiments and find little evidence of reputation-building behavior. Contrary to the expected result that reputation-building should lead to higher contributions when players' op-

---

<sup>5</sup>Selten and Stoecker use parameter estimates from the first 20 supergames to predict the outcomes of the last five supergames and find strong agreement between the predictions and actual outcomes.

<sup>6</sup>Subjects in Selten and Stoeker's experiment participated in 6-person matching groups so could have learned their opponents' types over 25 repetitions. Since subjects were told that they would play each opponent only once, however, they ruled this possibility out.

<sup>7</sup>In an fMRI study of a 10-period trust game—which is similar to the FRPD game—King-Casas et al (2005) find that second-movers' brains eventually signal the intent to cooperate before the first-movers' actions are revealed. They also become more accurate in predicting first-movers' actions. This is consistent with the hypothesis that players build a model of their opponent over time, though the data are not informative about the content of that model. A related idea is explored by Kahn and Murnighan (1993), who conduct an experiment on FRPDs in which they explicitly induce uncertainty about opponents' types by varying their pecuniary payoffs. They find that "weak" players (players for whom defection is not a dominant strategy in the stage game) are more cooperative than "strong" players (with typical prisoners' dilemma payoffs), and that uncertainty about opponents' payoffs increases cooperation for "weak" players.

Samuelson (1987) shows that cooperation can be sustained for at least some periods when the assumption that the number of periods is common knowledge is relaxed. Following this approach, Normann and Wallace (2012) experimentally compare repeated prisoners' dilemma games with known, random, and ambiguous number of periods, finding no significant differences in cooperation. An experiment by Bruttel, Güth, and Kamecke (2012) studies an FRPD in which the number of periods is uncertain. They find that cooperation breaks down closer to the final round than in a baseline treatment with a commonly-known finite horizon. They also find that many players cooperated after they were privately informed about the number of remaining periods. In the current study, the number of periods is publicly announced to all subjects to eliminate such uncertainty.

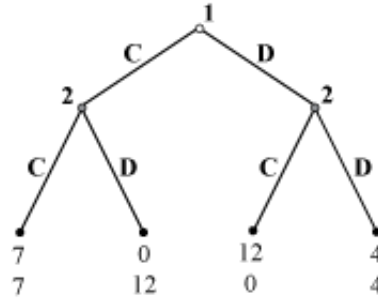


Fig. 1 A single stage of the sequential-move FRPD.

ponents are fixed, this is not always the case (Andreoni, 1988), and a restart effect seems more influential than the matching regime in raising contributions.

### 3 A REPUTATION-BASED THEORY OF COOPERATION

In this section we apply the theory developed by Kreps et al (1982) to the sequential-move FRPD played by the participants in our experiment and characterize the optimal strategies.<sup>8,9</sup> A single stage of the sequential-move FRPD played in our experiment is shown in Figure 1. As with Kreps et al (1982), we assume the first-mover believes the second-mover may be a tit-for-tat behavior type with positive probability and a rational second-mover is aware of (and can take advantage of) this belief.<sup>10</sup> The tit-for-tat type always reciprocates the first-mover's action, regardless of the period.<sup>11</sup> In the following analysis we therefore focus only on the rational, payoff maximizing second-mover's decisions.<sup>12</sup>

Let  $p_t$  be the first-mover's period- $t$  belief probability that the second-mover is the tit-for-tat type, with  $p_1 \in (0, 1)$  representing his prior belief. Updating occurs in equilibrium according to Bayes' rule. If  $p_t = 0$  in any  $t$ , then by standard unraveling arguments both players must play  $D$  in period  $t$  and thereafter. For this game there exist sequential equilibria in which there is a belief threshold  $\bar{p}_t$  such that the first-mover is willing to trust the second-mover (by playing  $C$ ) in period  $t$  if and only if  $p_t \geq \bar{p}_t$ . The rational second-mover prefers to maintain a reputation for being the tit-for-tat type by responding to  $C$  with  $C$  (i.e., conditionally cooperate) until some later round that is dependent on  $p_t$ . As the end of the game nears the expected payoff to the first-mover from continuing to play  $C$  declines since there are fewer rounds left and the probability that a rational second-mover plays  $D$  increases. In consequence, the belief threshold  $\bar{p}_t$  is strictly increasing in  $t$  up to  $\bar{p}_{10} = 4/7$ , given the specific stage-game payments shown in Figure 1.

When  $p_1 < \bar{p}_1$  the first-mover will defect in all periods, but when  $p_1 > \bar{p}_1$ , the first-mover initially trusts the second-over, who will perfectly imitate a tit-for-tap type if the is rational. Because

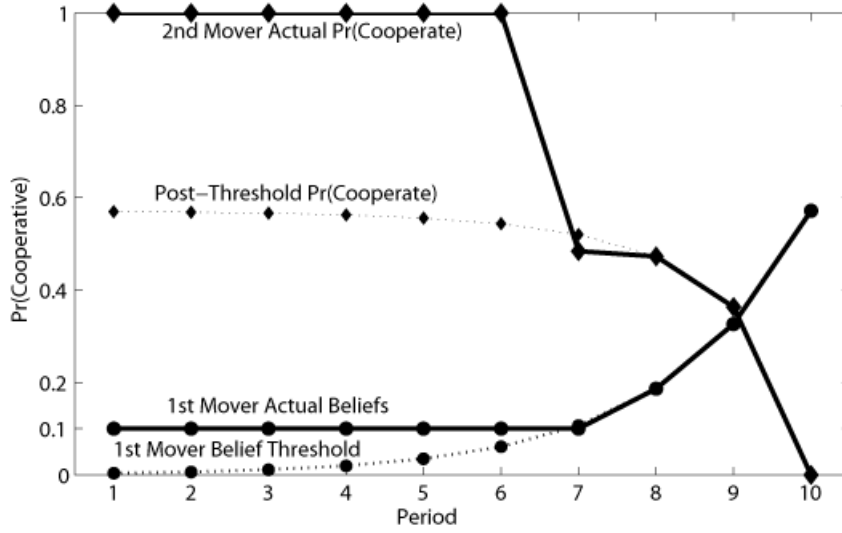
<sup>8</sup>See Kreps et al (1982) for a detailed derivation of the sequential equilibrium and see the appendix for a derivation of the equilibrium for the specific game used in this experiment, which includes the function used to generate Figure 2.

<sup>9</sup>We utilize a sequential move game so that we can focus on one player, the second-mover, and her ability to build a reputation by mimicking a behavioral type.

<sup>10</sup>That is, similar to Kreps et al (1982) the first-mover's prior is common knowledge. This assumption could be relaxed. For example, one could model the first-mover's prior as coming from a distribution of priors, in which case the second-mover's optimal strategy will be a function of the first-mover's expected prior. The distribution for this expectation may change as  $t$  increases too, but the intuition for the optimal strategies remains largely the same and revealing oneself as rational removes the uncertainty over the second-mover's type.

<sup>11</sup>Such beliefs are certainly justified given that tit-for-tat play is often observed in experimental data; see Andreoni and Miller (1993), for example.

<sup>12</sup>We will more simply refer to a player that has the objective of maximizing his own payoff as a rational player.



**Fig. 2** The sequential equilibrium of the sequential-move FRPD with a tit-for-tat type. See the appendix for a proof of the sequential equilibrium and the equations used to derive the figure.

the second-mover is either tit-for-tat, or imitating a tit-for-tat type, the first-mover's beliefs do not change in early periods. At some point in time, however, because the belief threshold increases as  $t$  increases,  $\bar{p}_t$  may rise above  $p_1$ , at which point the first-mover would stop trusting the second-mover. Let  $\bar{t}$  be the first period in which  $\bar{p}_t > p_1$ . The second-mover benefits from being trusted and would prefer to keep  $p_t$  weakly above  $\bar{p}_t$  in every period  $t \geq \bar{t}$ . He does this by playing a mixed strategy so that the first-mover's beliefs shift up to exactly  $p_t = \bar{p}_t$  for all  $t \geq \bar{t}$ . That is, if the second-mover is observed to conditionally cooperate at  $t$ , then the first-mover's belief that the second-mover is in fact a tit-for-tat type increases given that a rational second-mover would have defected with some probability. Formally, the second-mover conditionally cooperates in period  $t$  with probability

$$q_t^* = \frac{p_t}{1-p_t} \frac{1-\bar{p}_{t+1}}{\bar{p}_{t+1}},$$

and the first-mover's beliefs update to

$$p_{t+1} = \begin{cases} \frac{p_t}{p_t + q_t(1-p_t)} & \text{if C is realized;} \\ 0 & \text{otherwise.} \end{cases}$$

We refer to  $q_t^*$  as the post-threshold probability of cooperation. Since  $p_t$  must continue to increase over time,  $q_t^*$  must decrease accordingly to keep  $p_t$  below  $\bar{p}_t$ .

When  $p_t = \bar{p}_t$ , the first-mover is exactly indifferent between C and D. Since the second-mover is mixing, he must also be exactly indifferent between C and D. This is done by having the indifferent first-mover mix between C and D in period  $t+1$  with appropriate probabilities. If the first-mover's realized action is D then both types of second-mover respond with D and beliefs do not update. Thereafter the first-mover does not trust the second-mover (because  $p_{t+1} = p_t = \bar{p}_t < \bar{p}_{t+1}$ ), and the second-mover never has a chance to alter beliefs, so defection occurs in every subsequent

period.<sup>13</sup> The structure of a unique sequential equilibrium for a first-mover having belief  $p_1 > \bar{p}_1$  is displayed in Figure 2.

Observe that the path of equilibrium play depends crucially on  $p_1$ . If  $p_1 > \bar{p}_{10} = 4/7$  then no mixing phase is needed; both players play C with certainty until the rational second-mover defects in the final period. If  $p_1 \in [\bar{p}_1, \bar{p}_{10}]$  then mixing will begin in period  $\bar{t}$  (the smallest  $t$  such that  $p_1 \leq \bar{p}_t$ ), after which one defection will lead to defection in all subsequent actions. If  $p_1 < \bar{p}_1$  then the first-mover will never trust the second-mover, the second-mover will never have an opportunity to alter beliefs, and defection will occur in every action.

Regardless of  $p_1$ , the realized path of play must feature a regime shift from cooperation to defection. This shift can be triggered by either player, can occur in the first action, and may never occur if tit-for-tat players truly exist. In the laboratory beliefs are not directly observable and second-movers' types are unknown, so the time at which defection begins cannot be predicted without additional data.<sup>14</sup>

#### 4 EXPERIMENTAL DESIGN

Our experiment is designed to test directly the reputation-building aspect of the Kreps et al (1982) theory in a sequential-move, finitely-repeated prisoner's dilemma. In all treatments, 20 subjects are divided into two equal-sized groups: first-movers and second-movers. The sessions are divided into two Blocks. In each Block, each subject plays five finitely-repeated prisoners' dilemma supergames, with each supergame played against a different subject from the other group.<sup>15</sup> Through the course of the experiment, each subject will therefore play a supergame with each subject in the other group exactly once (5 in Block 1 and the other 5 in Block 2).<sup>16</sup> Each supergame is 10 rounds in length.

The sequential-move game allows us to focus on the second-mover and her opportunities for reputation building. We study this by varying the information structure in two treatments, denoted 2S and 1S. The first Block of five supergames is identical across the two treatments. Players see their opponent's history in the current supergame, but not from any prior supergames. In Block 2 of treatment 2S, subjects in each supergame see their opponents' entire history of play from their five Block 1 supergames, as well as the play of their opponents' opponents in these five supergames. Thus, all of the first-mover's actions, all of the first-mover's opponents' actions, all of the second-mover's actions, and all of the second-mover's opponents' actions from Block 1 are revealed to both players in each Block 2 supergame of 2S. It is commonly known that this information is revealed to both players.

In treatment 1S, only the first-mover's history from Block 1 is revealed to the second-mover (including the first-mover's actions and the first-mover's opponents' actions); the second-mover's history is not revealed to the first-mover. Again, this revelation structure is commonly known. The

<sup>13</sup>The second-mover cannot attempt to restore cooperation by playing C in response to D, for this would reveal his true rationality with certainty and result in defection in all subsequent periods.

<sup>14</sup>The model can be generalized slightly by allowing the first-mover's underlying beliefs to be non-stationary. The implications of the model are, however, similar to the stationary version.

<sup>15</sup>Having subjects play multiple supergames in each block allows them to become familiar with the game, and more importantly, allows second-movers to reveal more information about their types in Block 1. In any single supergame, it is possible that an individual second-mover might face a very uncooperative first-mover, so that no information about the second-mover's type would be revealed.

<sup>16</sup>In a couple of sessions, less than 20 subjects participated. In these sessions, subjects played against a different subject from the other group until they had played against all of them once. In the remaining supergames of Block 2, subjects were matched randomly with one of the subjects had already faced in Block 1 but not yet faced in Block 2.



purpose of this treatment is to control for changes in behavior between Blocks 1 and 2 that result from revealing the first-mover's Block 1 history of play to Block 2 second-movers as well as to use as a control for any "restart effect". Differences between 1S and 2S can then be interpreted as the effect of revealing the second-mover's Block 1 history of play to Block 2 first-movers (revealing reputations), given that the second-mover is equally informed about the first-mover's history. The treatment names 2S and 1S are mnemonic for "two-sided" and "one-sided" knowledge of histories, respectively.

At the beginning of Block 1 of both treatments, subjects know they will play five finitely-repeated prisoners' dilemma supergames against five different opponents in Block 1. Subjects are also told at the beginning of Block 1 of both treatments that we will conduct a second experiment immediately following the first in which they may also participate if they want. They are informed that instructions for the second experiment will be distributed after the first experiment concludes. The subjects are not told that they will play prisoners' dilemmas in the second experiment or that their histories from the Block 1 may be revealed in Block 2. Only at the beginning of the Block 2 do they learn that they will play additional repeated prisoners' dilemma supergames and that their history from Block 1 will be revealed (depending on the treatment). Subjects have the option of taking their earnings after Block 1 and leaving instead of participating in the Block 2 experiment, and if they participate in the second experiment they are guaranteed a second minimum show-up payment. All of the subjects chose to stay for the second experiment.

As mentioned in the previous section, we use the same stage-game payoffs as Andreoni and Miller (1993), shown in Figure 1. Subjects are paid for one randomly-selected supergame out of the five supergames in each Block, and this is known in advance.

We use a strategy-elicitation method for second-mover choices, which asks subjects to enter their action conditional on the first-mover defecting (choosing right) and their action conditional on the first-mover cooperating (choosing left) before learning the first-mover's action in each round. The second-mover's strategy for a particular round is implemented for them after the first-mover chooses an action for that round. In a survey paper comparing the strategy method to direct-response, Brandts and Charness (2011) found that while the two methods can induce different behavior in some experiments, there is generally no significant difference in behavior in prisoner's dilemma experiments between the two. They conclude from their survey that in cases where behavior differs between the two methods, the strategy method provides a lower bound for treatment effects.<sup>17</sup>

In the sequential-move prisoner's dilemma, the second-mover's strategic intention may be censored in rounds where the first-mover defects, which leads to defection by the second-mover as well according to a wide range of strategies. We elicited second-mover choices using the strategy method so that we would be able to identify the second-mover's intent to cooperate if reciprocated,

<sup>17</sup>Experiments have also been conducted on the (one-shot) sequential prisoners' dilemma and they generally show little difference from simultaneous-move setups. Bolle and Ockenfels (1990) found little difference in cooperation levels between simultaneous and sequential one-shot prisoners' dilemma using the strategy method to elicit second-mover strategies. Brandts and Charness (2000) found no significant difference in cooperation between the sequential one-shot prisoner's dilemma using the strategy method and direct response. Blanco, Engelmann, Koch, and Normann (2011) used the strategy method with role uncertainty in several information conditions and belief-elicitation treatments to show that correlation between strategies in different roles is driven partially (but not completely) by a consensus effect. Clark and Sefton (2001) examined sequential prisoners' dilemma games with varying levels of temptation and overall stakes in both the United States and the United Kingdom. They found substantial cooperation levels in early rounds which diminished by the tenth and final round. They also found that second-movers were much more likely to cooperate if the first-mover cooperated, but this tendency also decreased across rounds and with higher temptation levels. Higher overall stakes lead to a slight increase in second-mover reciprocal cooperation in the UK, but a decrease in reciprocal cooperation in the US.

**Table 1** Summary Information

Treatment	Revealed Information	Number of Sessions	Number of Subjects
1S	History of play for the first-mover and his opponents from block 1.	7	130
2S	History of play for both the first and second-movers and their opponents from block 1.	7	134

regardless of the first-mover’s actual choice. This way, if the first-mover defected and the second-mover responded by defecting in the same round, we would know whether the second-mover would have continued cooperating or unilaterally defected if the first-mover had instead cooperated in that round. All of the second-mover results reported in the paper focus on this conditional cooperation (the second-mover’s choice conditional on cooperation by the first-mover) and not the choice actually implemented for the second-mover in response to the first-mover.

## 5 RESULTS

We conducted 14 sessions of the experiment at The Ohio State University experimental economics lab with a total of 264 subjects: seven sessions of the 1S treatment with 130 subjects and seven of the 2S treatment with 134 subjects. Subjects were chosen randomly from a pool of students at The Ohio State University who had previously signed-up to be considered for participation in economic experiments. Subjects could not participate in more than one session of this experiment. Table 1 provides summary information for the experiment. Payoffs in the experiment were denominated in “points” and converted into dollars at the rate of 4 points per dollar. Average earnings per subject were approximately \$27. We proceed by first analyzing the players’ aggregate behavior and then focus on how information impacts behavior at the individual level.

### 5.1 Aggregate Behavior

We focus on cooperation by the first-mover and conditional cooperation by the second-mover (choosing to cooperate conditional on the first-mover cooperating) as our outcomes of interest. The first hypothesis, based on Kreps et al (1982), states that cooperation will be eliminated when the second-mover’s history reveals she is rational.

**Hypothesis 1** *Compared to Block 2 of 1S, the rate of aggregate cooperation should not be higher in Block 2 of 2S.*

Figure 3 displays the paths of aggregate cooperation by first-movers over the course of a supergame. The first row represents the data pooled from each treatment for the first block, the second row represents the Block 2 data for 1S, and the third row represents the Block 2 data for the 2S treatment. Each column represents a supergame in the order in which they were played. The paths of play in Block 1 (1S and 2S pooled) and Block 2 of 1S are quite similar for first-movers as we would expect given that there is no change in the information given to first-movers between these treatments. The path of play in Block 2 of 2S, however, exhibits higher levels of cooperation than the other treatments. Moreover, the level of cooperation increases across supergames and the path of play shows a clear endgame effect in late supergames of Block 2—something that is

not clearly apparent in other supergames where aggregate cooperation declines more steadily by round.

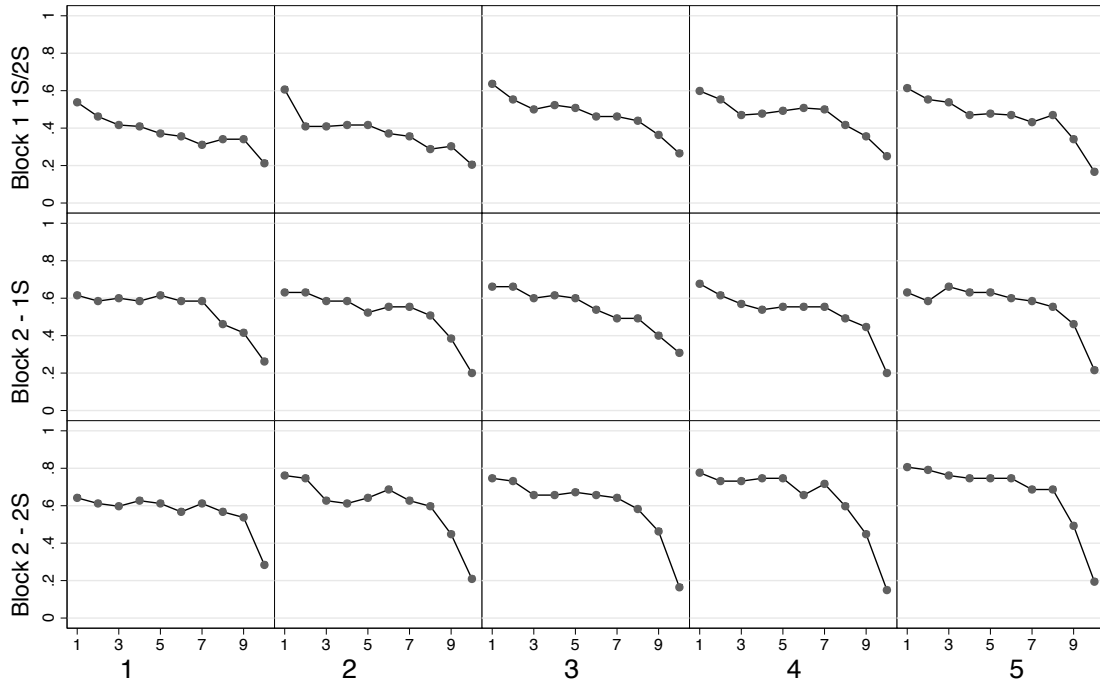


Fig. 3 first-mover average cooperation.

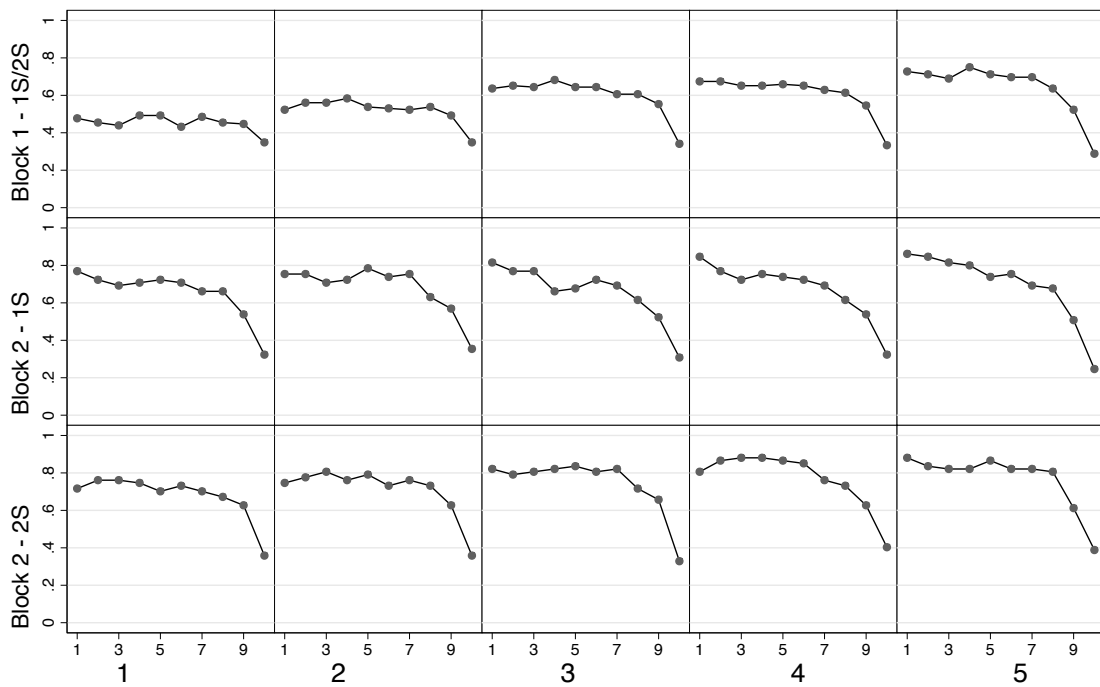


Fig. 4 second-mover average conditional cooperation.

**Table 2** Aggregate Cooperation

	1st Movers		2nd Movers	
	# of Subjects	% Cooperation	# of Subjects	% Conditional Cooperation
Block 1				
1S	65	43.2%	65	56.4%
2S	67	42.9%	67	57.0%
Block 2				
1S	65	53.4%	65	67.5%
2S	67	60.9%	67	73.1%

Figure 4 displays the paths of aggregate second-mover conditional cooperation—cooperate conditional on the first-mover cooperating—over the course of a supergame. The paths of play are again shown by treatment and supergame, as in Figure 3. The paths of aggregate cooperation over supergames in Block 2 of 1S show more of an endgame effect than those in Block 1 (pooled) since cooperation is maintained in early rounds then drops off more sharply in later rounds. However, round 10 behavior is roughly the same across blocks, suggesting that players view the endgame similarly in both.<sup>18</sup>

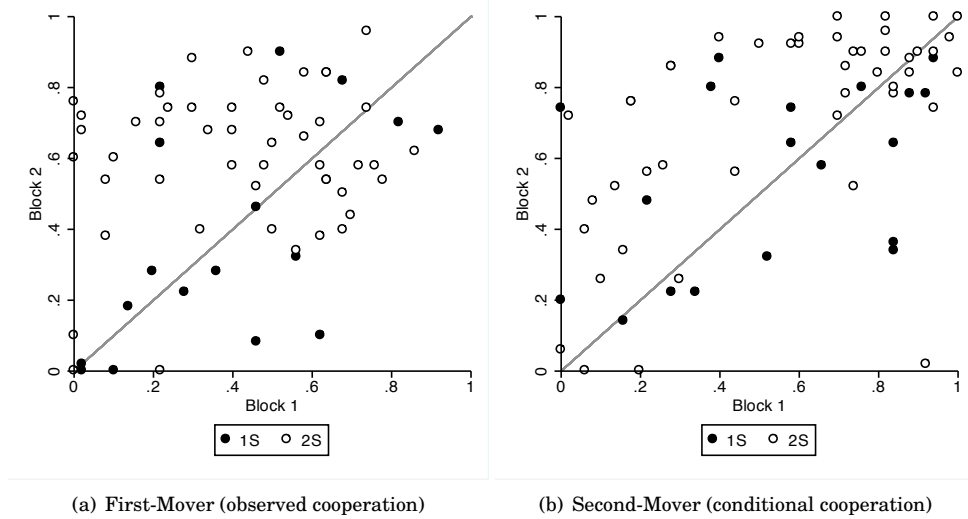
The paths of play in Block 2 of both 1S and 2S, exhibit more cooperation than occurs in Block 1 with 2S exhibiting even higher levels of cooperation than that in Block 2 of 1S despite first-movers observing the histories of second-movers. That greater overall cooperation and sustained cooperation until later rounds of the supergames are observed when second-mover histories are revealed to first-movers is again inconsistent with the predictions of Kreps *et al* (1982).

**Result 1** *We observe greater aggregate cooperation in Block 2 when second-movers' histories of play are exposed to first-movers (treatment 2S, compared to 1S).*

Table 2 reports the aggregate cooperation frequencies of first-movers and the aggregate conditional cooperation frequencies of second-movers by block and treatment. There is almost no difference in overall cooperation between Block 1 cooperation in 1S and 2S for either first and second-movers. However, not only do both first- and second-movers cooperate more in Block 2 of 2S, but their cooperation rates are slightly higher than that of first- and second-movers in Block 2 of 1S. Wilcoxon-Mann-Whitney tests (bootstrapped to account for clustering by session) confirm that there is no between-treatment difference in Block 1 cooperation rates for first- or second-movers (p-values 0.959 and 0.662, respectively), while there is slightly higher Block 2 cooperation rates in 2S than in 1S (p-values are 0.063 and 0.071 for first- and second-movers, respectively).<sup>19</sup>

<sup>18</sup>For both first and second-movers in both treatments, round 10 cooperation/conditional-cooperation is never significantly higher in block 2 than in block 1 (Wilcoxon signed-rank tests with the unit of observation being the subject-level average cooperation/conditional-cooperation in round 10 across all supergames in a block).

<sup>19</sup>The unit of observation in these tests is the average (conditional) cooperation for an individual first (second) mover across all periods of all Block 2 supergames. As we are testing whether there is higher cooperation in 2S than 1S, the latter p-values are based on a 1-sided test in which the null is  $\mu_{2S} = \mu_{1S}$  and the alternative is that  $\mu_{2S} > \mu_{1S}$ . The two-sided p-values are 0.126 and 0.141, respectively.



**Fig. 5** Average cooperation by mover by subject (all supergames).

## 5.2 Individual-Level Behavior

As Cooper et al (1996) point out, individual behavior provides a better test of theoretical predictions in the FRPD than aggregate cooperation. We now analyze individual-level behavior, and classify players into types based on that behavior.

Figures 5(a) and 5(b) plot average cooperation for each subject in Blocks 1 and 2. The first figure shows cooperate rates for first-movers while the second shows conditional cooperation rates for second-movers. Each data point represents a single subject. Subjects in the 1S treatment are black circles, while subjects in the 2S treatment are white. Lack of a treatment effect would express itself in data points scattered symmetrically about the 45-degree line. In both figures the 1S data is distributed in this manner. However, the 2S data is skewed above the 45-degree line, indicating that revealing Block 1 histories of second-movers increases cooperation for both first and second-movers. Wilcoxon signed-rank tests for data having within-group correlations (Larocque, 2005) show that differences in frequency of cooperation across blocks are highly significant in 2S ( $p$ -value=0.063 for first-movers and  $p$ -value=0.015 for second-movers) and insignificant in 1S ( $p$ -value=0.222 for first-movers and  $p$ -value=0.210 for second-movers).<sup>20</sup>

We classify players into types within each block to control for the heterogeneity in individual behavior. By focusing on how information differentially impacts players conditional on their type, we can identify whether the behavior of some players is consistent with the reputation-building theory. Importantly, we classify players based on histories of play only—not on the off-path choices of second-movers—as our main focus is on the information contained in Block 1 histories revealed in Block 2. Furthermore, we classify second-movers based solely on their behavior in rounds in which the first-mover cooperated. Second-movers cannot reveal any information about their types in rounds in which the first-mover defected, as tit-for-tat and reputation-building players alike would defect in such cases. For ease of explaining the type classification, rounds in which the first-mover cooperated are referred to as *trusting* rounds.

<sup>20</sup>We perform a power calculation for the 1S test, assuming the effect size found in 2S and a 10% significance level, and find that our power is 85%. The standard threshold for acceptable power is 80%.

**Table 3** Second Mover Type Transition Matrix

1S		Block 2			
		Defector	Cooperator	Imitator	Total
Block 1	Defector	10	12	3	25 (38.5%)
	Cooperator	3	16	9	28 (43.1%)
	Imitator	1	3	8	12 (18.5%)
Total		14 (26.5%)	31 (47.7%)	20 (30.8%)	65 (100.0%)

2S		Block 2			
		Defector	Cooperator	Imitator	Total
Block 1	Defector	10	12	4	26 (38.8%)
	Cooperator	2	17	10	29 (43.3%)
	Imitator	0	9	3	12 (17.9%)
Total		12 (17.9%)	38 (56.7%)	17 (25.4%)	67 (100.0%)

The classification procedure works as follows. For each supergame, a second-mover is classified as an *Imitator* if her behavior is consistent with the reputation-building strategy of Kreps et al (1982); i.e., she plays cooperate in the first trusting round and continues doing so until some later trusting round (possibly round 10), after which she plays defect in each subsequent trusting round. Otherwise, she is classified as a *Cooperator* if she cooperates in the first trusting round (but did not play as an *Imitator*) or as a *Defector* if she does not cooperate the first trusting round. Her type classification for the block is identified by the mode of her five supergame classifications within that block. In the case of a tie, the most recently-used modal type is used.<sup>21</sup>

By this procedure we arrive at a classification of second-movers that summarizes quite well their Block 1 history of play as observed by first-movers. The overall percentage of Block 1 supergames having a second-mover with a type classification of *Imitator*, *Cooperator*, or *Defector* is 18.2%, 41.7%, and 39.4%, respectively.<sup>22</sup> Based on these type classifications we have the following hypothesis regarding how a second-mover's classification will change in Block 2 when her history of play is revealed.

**Hypothesis 2** *Second-movers in 2S who play as Imitators or Defectors in Block 1 play as Defectors in Block 2 because they are revealed to first-movers as rational.*

Table 3 reports the type classifications for second-movers by treatment and block in the form of a Markov transition matrix. The data reveal that transitions between second-mover types between Blocks 1 and 2 in 2S are inconsistent with reputation-building. Surprisingly, nearly half of the Defectors in Block 1 become Cooperators in Block 2 and only 38% of the Block 1 Defectors remain Defectors in Block 2. Even more striking, none of the second-movers who are Imitators in Block 1 of 2S become Defectors in Block 2. These findings are summarized in the following result.

<sup>21</sup>In the appendix we present similar results in which players are classified only by their type from the last supergame they play.

<sup>22</sup>The proportion of second-movers who are classified as Cooperators may be inflated relative to Imitators because first-movers defected first in 41.5% of the Block 1 games with Cooperators. It is possible that these second-movers were using a reputation-building strategy that is not revealed because of the first-movers' defection; however, classifying them as Cooperators is still useful because second-movers are not revealed to first-movers as rational.

**Table 4** First Mover Type Transition Matrix

1S		Block 2		
		Non-Trusting	Trusting	Total
Block 1	Non-Trusting	21	10	31 (47.7%)
	Trusting	2	32	34 (52.3%)
	Total	23 (35.4%)	42 (64.6%)	65 (100.0%)
2S		Block 2		
		Non-Trusting	Trusting	Total
Block 1	Non-Trusting	7	19	26 (38.8%)
	Trusting	2	39	41 (61.2%)
	Total	9 (13.4%)	58 (86.6%)	67 (100.0%)

**Result 2** Only 26.3% of second-movers who are classified as Defectors or Imitators in Block 1 of 2S are classified as Defectors in Block 2, while 55.2% of them are classified as Cooperators in Block 2.

The first-mover type classification is slightly different. Kreps et al (1982) allows for two possible types of first-movers: those who believe the probability of an irrational second-mover is high enough to justify cooperation in round 1 of a supergame, and those who do not. Therefore, we classify a first-mover as *Trusting* if her modal behavior is to cooperate in round 1 of the five supergames in a given block. Otherwise, a first-mover is classified as *Non-Trusting*. Based on these type classifications we have the following hypothesis regarding how a first-mover's classification will change in Block 2 when the second-mover's history of play is revealed.

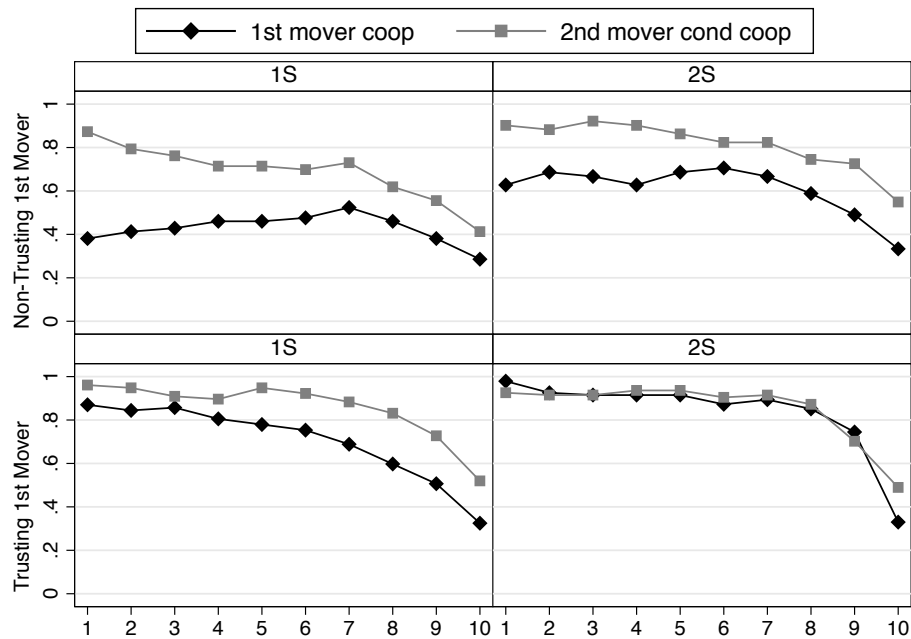
**Hypothesis 3** Compared to Block 2 of 1S, first-movers are less likely to be classified as *Trusting* in Block 2 of 2S.

Table 4 reports type classifications for first-movers by treatment and block in the form of a Markov transition matrix. In 1S, 47.7% of first-movers are Non-Trusting and about a third (32.3%) of these first-movers become Trusting in Block 2. However, in 2S more than 2/3 (73.0%) of first-movers who are Non-Trusting in Block 1 (38.8% of all first-movers in this treatment) transition to Trusting in Block 2. Because this result is not conditional on the revealed type of second-mover opponents, it would not contradict the theoretical predictions if no second movers were exposed as rational. But some were, so we should expect an aggregate decrease in the number of Trusting first movers. Instead, we find the opposite: First movers become more trusting when histories are revealed.

**Result 3** A higher proportion of first-movers are *Trusting* in Block 2 of 2S compared to Block 2 of 1S ( $p=0.044$ ), while there is no difference in the proportion of *Trusting* first-movers in Block 1 ( $p=0.463$ );<sup>23</sup> Additionally, there is a larger increase in the proportion of first-movers who are *Trusting* from Block 1 to Block 2 in 2S: 25.4 percentage point increase in 2S ( $p=0.032$ ) and a 12.3 ( $p=0.122$ ) percentage point increase in 1S.<sup>24</sup>

<sup>23</sup>P-values are based on a Wilcoxon-Mann-Whitney tests with cluster bootstrapped variances.

<sup>24</sup>P-values are from a McNemar test for data having within-group correlations (Durkalski et al, 2003). We calculate the power of the 1S test to be 78%, assuming the effect size observed in 2S.



**Fig. 6** Block 2 cooperation rates when the second mover is revealed to have been a Cooperator type in Block 1.

We now study first-mover reactions to the second-mover's revealed type. Because first-movers should only cooperate if there is some positive probability that the second-mover is irrational, exposing the second-mover as an Imitator or Defector in Block 1 should convince the first-mover that the second-mover is rational and destroy cooperation in Block 2 of 2S. This is summarized in the following hypothesis.

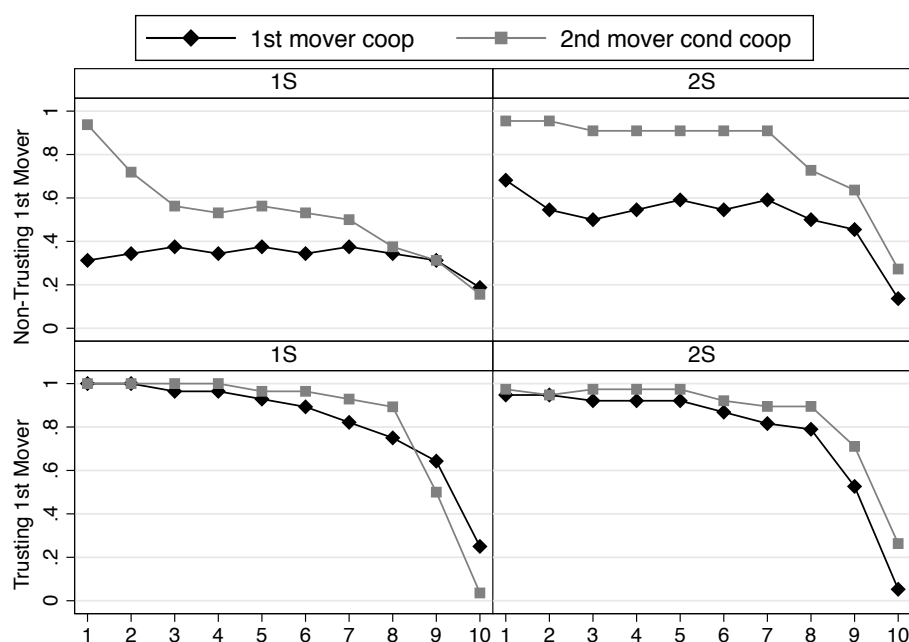
**Hypothesis 4** *In Block 2 of 2S, a first-mover whose opponent's Block 1 history is an Imitator- or Defector-type will play as a Non-Trusting type against this opponent.*

We test hypothesis 4 by splitting the Block 2 data by the second-mover's Block 1 type and examining how first-movers respond to their opponents' revealed histories. Note that in the following analysis we study Block 2 data using type classifications based on Block 1 data.

First we look at the case where the second mover's history indicates they are a Cooperator type. Figure 6 shows average first-mover cooperation and second-mover conditional cooperation by round in Block 2 games for only those supergames where the second-mover's Block 1 histories classify her as a Cooperator. The top two panels show the Block 2 cooperation rates in 1S and 2S, respectively, for supergames in which the first-mover is Non-Trusting. The bottom panels are for Block 2 supergames in which the first-mover is Trusting. Both Non-Trusting and Trusting first-movers cooperate more frequently and into later rounds after the second-mover is revealed to be a Cooperator compared to when no histories are revealed to first-movers.

Figure 7 presents a similar view of the data for Block 2 supergames where the second-movers' Block 1 histories classify her as an Imitator. First-movers with Trusting Block 1 histories have very similar cooperation rates in each corresponding round across treatments. Block 2 cooperation is much higher in 2S than in 1S for Non-Trusting first-movers, however. Revealing information that indicates that second-movers are willing to cooperate increases first-mover cooperation even if that information reveals that the second-mover is rational.



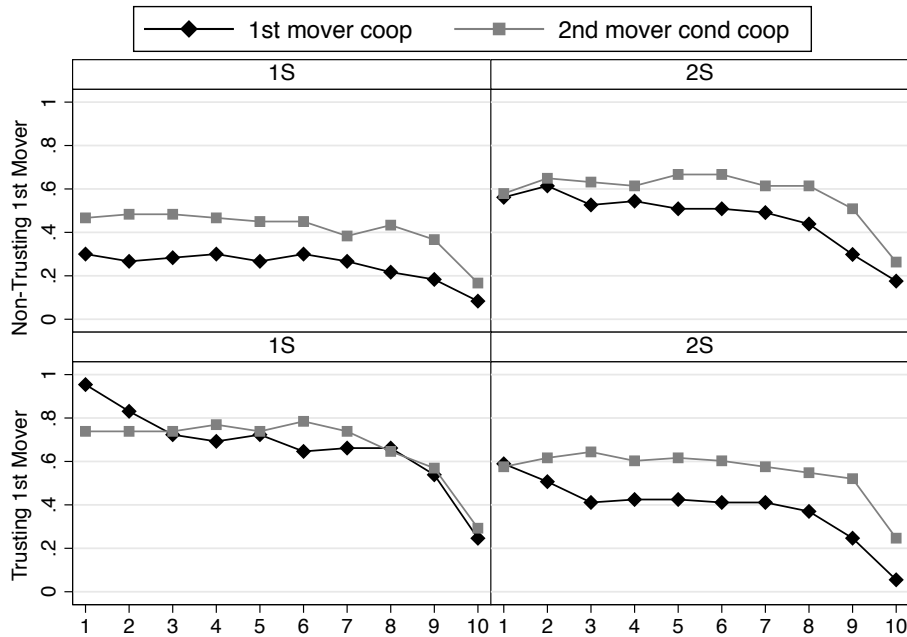


**Fig. 7** Block 2 cooperation rates when the second mover is revealed to have been an Imitator type in Block 1.

Figure 8 is similar to Figures 6 and 7 except data is from Block 2 supergames where the second-movers' Block 1 histories classify her as a Defector. Initial cooperation by Trusting first-movers in 2S decreases substantially between 1S and 2S when the second-mover is revealed to be a Defector. However, Non-Trusting first-movers cooperate slightly more when the second-mover is revealed to be a Defector. These findings are consistent with the notion that providing the second-mover's history causes first-movers to update their beliefs regarding the degree of cooperation that they can expect from a first-mover. Gächter and Thoni (2005) find a similar result in their public goods game experiment in which players' behavior from the first game is revealed to other players in the subsequent series of games (and players were not told this would happen before the first game). They find that when players whose contributions were low in the first game are grouped with other low contributors in subsequent games, they contribute more than when grouped randomly. As in our experiment, this behavior may be due to a failure of backwards-induction, combined with updating of pessimistic beliefs that these low contributors hold in the first game; an interpretation consistent with the model of naive beliefs that we develop in Section 6.

### 5.3 Rate of First Defection

We use logit regressions to test how first-mover cooperation in Block 2 supergames depends on the type of second-mover. Robust standard errors are corrected to account for within-subject and within-session correlations between observations; i.e., the standard errors are clustered by individual first-mover and session. Specifications (i) and (ii) of Table 5 show the impact of revealed player histories on first-mover initial cooperation. Specification (i) explores the impact of revealing a player's history of play and indicates that first-movers are more trusting when second-movers are observed to have cooperated in block 1 supergames. Specification (ii) examines the impact that a second-mover's type has on the likelihood that a first-mover will be Trusting. In this speci-



**Fig. 8** Block 2 cooperation rates when the second mover is revealed to have been a Defector type in Block 1.

fication, only second-mover types from 2S are identified by first-movers while second-mover types from 1S are not observed and represent the omitted second-mover type in the regression.<sup>25</sup> The estimates show that both Trusting and Non-Trusting first-movers are more likely to trust in Block 2 when the second-mover is revealed to be a Cooperator or an Imitator, compared to when no history is revealed. However, when facing a Defector, Non-Trusting first-movers are no more likely to trust than when no history is revealed. Regardless of the second-mover's type, Trusting first-movers are more likely than Non-Trusting first-movers to trust in Block 2 when the second-mover's history of play is revealed. In this case, even though a second-mover's model type is Defector, the first-mover may be responding to the fact that the second-mover likely cooperated in some supergames.

Specifications (iii)-(vi) of Table 5 explore the impact of revealing player types on the number of periods of joint cooperation before first defection in Block 2 supergames. As the number of rounds of initial cooperation are censored between 0 and 10, the results are estimated using Tobit regressions. Furthermore, in contrast to specifications (i) and (ii), second-movers are classified based on their Block 1 history of play in *both* 1S and 2S. Although the first-mover cannot classify a second-mover in 1S, it is useful as an explanatory variable as a second-mover's type will impact the overall degree of cooperation and will help us to assess what player characteristics generated the cooperation rates.

Regardless of whether the second-mover's history is revealed, average cooperation is slightly longer in block 2, with some weak evidence that it is even longer in 2S, though the estimate is not statistically different from zero at conventional levels. That cooperation is longer in block 2 of 1S suggests that some of the increased cooperation rates seen in 2S is likely driven by second-movers conditionally cooperating for longer than in block 1. Though average cooperation is informative, it can be biased by the player types. If, for example, one treatment included more cooperative

<sup>25</sup>As the type effects are relative to second-mover types in 1S the estimates do not simply reflect a restart effect.

**Table 5** The Effect of Player Types on First-Round Cooperation and Rounds of Cooperation before First Defection

	1st Round Cooperation by First-Mover		# of Rounds of Cooperation Before 1st Defection			
	(i)	(ii)	(iii)	(iv)	(v)	(vi)
Avg. Block 1 Cooperation	0.054 (0.308)		0.538 (0.021)			
× 2S	0.138 (0.019)		0.322 (0.154)			
Defector		-0.188 (0.651)				
Imitator		0.937 (0.046)		1.606 (0.186)	2.116 (0.012)	1.586 (0.344)
Cooperator		1.069 (0.040)		3.118 (0.004)	0.274 (0.818)	3.722 (0.011)
Trusting						
× Defector		1.662 (0.001)		9.779 ( $< 0.001$ )		
× Imitator		3.450 ( $< 0.001$ )		10.472 ( $< 0.001$ )		
× Cooperator		2.531 ( $< 0.001$ )		7.002 ( $< 0.001$ )		
2S						
× Defector				4.647 ( $< 0.001$ )	-4.888 ( $< 0.001$ )	5.863 ( $< 0.001$ )
× Imitator				3.997 (0.072)	-1.360 (0.011)	5.628 (0.055)
× Cooperator				3.980 (0.074)	1.996 (0.232)	5.334 (0.071)
2S×Trusting						
× Defector				-10.164 ( $< 0.001$ )		
× Imitator				-5.744 (0.010)		
× Cooperator				-1.855 (0.214)		
Supergame#	0.385 (0.036)	0.532 (0.023)	0.755 (0.122)	0.633 (0.243)	-0.566 (0.227)	3.433 (0.004)
(Supergame#) <sup>2</sup>	-0.047 (0.122)	-0.066 (0.092)	-0.082 (0.286)	-0.065 (0.439)	0.098 (0.201)	-0.440 (0.019)
Constant	-0.368 (0.241)	-1.257 (0.001)	-0.563 (0.506)	-3.916 ( $< 0.001$ )	7.510 ( $< 0.001$ )	-10.574 ( $< 0.001$ )
Observations	660	660	660	660	375	285

Notes: Column (v) includes observations with Trusting first-mover and column (vi) includes observations with Non-Trusting first-mover only. p-values from robust standard errors clustered by individual first-mover and session in parentheses.

types, then this could make it appear that average cooperation increased more for that treatment. Specifications (iv)-(vi) are designed to explore how first- and second-mover types relate to the duration of cooperation. Specification (iv) includes all of the data while specification (v) includes only Trusting first-movers and specification (vi) includes only Non-Trusting first-movers to help clarify the impact of the first-movers' types. The estimate for Cooperator in specification (iv) reveal that playing with a more cooperative second-mover in either 1S or 2S causes a first-mover to cooperate for longer within a supergame, especially when the first-mover is Non-Trusting.<sup>26</sup> Interestingly Non-Trusting first-movers are also more likely to cooperate longer with Defector and Imitator second-movers. Given the results in (ii), part of this result is likely driven by Non-Trusting first-

<sup>26</sup>For Trusting first-movers the difference in cooperation between 2S and 1S is 2.125, but only with  $p = 0.259$ .

movers' trusting second-movers more after observing that these second-movers did not initially defect. Trusting first-movers, however, do not cooperate as long in 2S when facing a Defector and, to a lesser extent, an Imitator type second-mover. As these first-movers have shown that they are more willing to trust, even when the second-mover's history of play is not known, this result can likely be explained by first-movers examining the number of rounds that second-movers cooperated and trying to pre-empt their defection without performing full backward induction.<sup>27</sup>

Specifications (v) and (vi) help to clarify the impact of the first-mover's type by restricting the sample to either Trusting or Non-Trusting first-movers. Reinforcing the results in (iv), the estimates for specification (v) indicate that first-movers update the amount that they are willing to cooperate based on the amount second-movers cooperated in Block 1. Recall that defector is the omitted category, thus when trusting first-movers encounter a Defector or Imitator in Block 2 they will cooperate for fewer rounds than in 1S but they will also cooperate more with either an Imitator or Cooperator than a Defector in 2S. The results of specification (vi) reveal that Non-Trusting first-movers, who were more pessimistic in their Block 1 beliefs (i.e., have a high prior that their opponent is not a tit-for-tat type), cooperate more in Block 2 even when facing a Defector who has revealed that she is not a tit-for-tat type, suggesting that first-movers update their beliefs about how long these second-movers will cooperate based on the history of play. While the amount of additional cooperation for a Non-Trusting first-mover does not depend on second-mover type, revealing the second-mover's history of play is sufficient to increase average cooperation with a Non-Trusting first-mover.

In summary, both Trusting and Non-Trusting first-movers are more trusting when they can view the second-mover's history of play; Non-Trusting first-movers cooperate longer in 2S than in 1S regardless of the second-mover type; and Trusting first-movers cooperate for fewer periods with Defectors and to a lesser extent with Imitators. These results indicate that cooperation in the finitely-repeated prisoners' dilemma continues even after the opportunity for reputation-building by second-movers is destroyed. We therefore reject Hypothesis 4.

**Result 4** *Compared to a first-mover whose opponent's history is not revealed, a first-mover whose opponent's Block 1 history is revealed as Imitator-type (a) is more likely to cooperate in round 1 of a Block 2 supergame, and (b) Non-trusting first-movers cooperate longer when second-mover histories are revealed while Trusting first-movers do not cooperate for as long when second-movers are revealed to be Defectors or Imitators.*

#### 5.4 Beliefs

In some later sessions, we added a belief elicitation stage at the end of the experiment. Subjects were presented with a sequence of 5 randomly-selected Block 1 histories of first and second movers from previous sessions of the same treatment. They were then asked to state how many rounds they believed these past participants cooperated before the first defection. For one randomly-selected belief elicitation question, each subject was paid \$5 if her stated belief was exactly correct.<sup>28</sup>

<sup>27</sup>This explanation is also consistent with the rates of defecting first shown in Table A.1: in 2S, Trusting first movers facing Imitator-type second movers defect first 68.4% of the time, compared to 42.0% in 1S.

<sup>28</sup>We decided not to elicit beliefs before or during gameplay because doing so is difficult and may itself affect beliefs and behavior. We thank an anonymous referee for the suggestion of eliciting beliefs about third parties.

Regression results are shown in Table 7 in the Appendix, summarizing how elicited beliefs responded to observed cooperation rates in displayed histories as well as the role of the player whose beliefs are elicited. While we find some evidence that elicited beliefs respond to observed cooperation in the expected directions, statistical significance is weak due to small sample size. Furthermore, we find that first and second movers report systematically different beliefs, suggesting that elicited beliefs may be biased by the subjects' own experience in the experiment.

## 6 A MODEL OF NAÏVE BELIEFS

In contrast to the predictions of the reputation-building theory of Kreps et al (1982), our experiment shows that there will be substantial cooperation in FRPDs even when players' previous histories are revealed. Moreover, learning that an opponent imitated a tit-for-tat type frequently increases cooperation, rather than destroying it.

The observed increases in cooperation may be partially motivated by fairness or indirect reciprocity. For example, Ho and Su (2009) utilize an ultimatum game in an experiment and find that roughly half of their subjects value how their offer compares to offers other followers have received. A similar notion of fairness could plausibly affect cooperation in other games, such as the FRPDs played here, when subjects can view how their current opponents treated previous opponents; however, fairness would seem to be a less salient concern in a prisoner's dilemma setting, where payoffs are determined more symmetrically (i.e., both players have a hand in the outcome), than in an ultimatum game setting, where control over payoffs is heavily unbalanced. Moreover, a concern for fairness would not explain why Non-Trusting first-movers become more Trusting when they observe that their opponent is a Defector. In community games with public reputations, indirect reciprocity has been shown to induce cooperation (see Nowak and Sigmund, 2005, for a recent survey), but these games typically feature frequent re-matching among small groups of subjects, who play simple, one-shot stage games with unilaterally determined payoffs, making reputation in future matches a dominant concern. In contrast, the matching protocol we employ minimizes the opportunity for indirect reciprocity by generally matching each subject no more than once with each other subject. In addition, the payoff gradient of strategic interactions in an FRPD supergame is large enough that the direct payoff consequences of behavior in a given supergame should reduce the extent to which behavior is motivated by potential payoffs in future supergames.

Though we cannot rule out that cooperation in our experiment may be partially motivated by fairness or indirect reciprocity, we have reason to doubt that they are the primary drivers of the cooperation we observe. For example, our results are consistent with the findings of Kagel and McGee (2014)'s experiment on FRPDs played by teams. They observe an increase in cooperation over time because teams anticipate that their opponent will defect in later rounds. Kagel and McGee find no evidence in team chat logs, however, that teams anticipate opponents' beliefs and attempt to mimic an irrational player as the Kreps et al (1982) model predicts. Nor do they find evidence of reciprocity or concerns for fairness. Instead, they find evidence that teams anticipate when their opponents will defect (i.e., their opponents' actions, not their beliefs) and attempt to defect one round earlier, which they interpret as a failure of common knowledge of rationality. These results suggest that boundedly rational behavior is a primary driver of cooperation in FRPDs. We

now propose a simple model of boundedly rational behavior, in which we relax the requirement that players' prior beliefs be consistent with an opponent's best response.<sup>29</sup>

As in Kreps et al (1982), players in our model decide in which round to stop playing tit-for-tat and begin unconditionally defecting. This is decided by weighing the long-term benefit of cooperation against the risk of the other player defecting first.<sup>30</sup> The difference between this approach and that of Kreps et al (1982) is that players do not engage in higher level reflection about the beliefs of their opponent. Instead, players form "naïve" beliefs which may not be consistent with their opponent's best response. This is consistent with Kagel and McGee (2014), who provide evidence that players simply try to defect one round before their opponent's anticipated first defection without reflecting on the beliefs of their opponents. Given arbitrary initial beliefs about how many rounds the opponent will continue playing tit-for-tat, players update beliefs within the game based on their opponents' choices using Bayes' rule and choose the optimal round to stop playing-tit-for-tat and begin unconditionally defecting.<sup>31</sup>

This simple model generates predictions that are consistent with our experimental results. First, the common finding that cooperation in the FRPD does not break down until one or two rounds before the last is consistent with the model if players have uniform or even more pessimistic beliefs.<sup>32</sup> Second, the first-mover who plays cooperatively in Block 1 may defect earlier in Block 2 supergames after the second-mover's history of play from Block 1 exposes him as generally uncooperative. Third, there will be as much, if not more, cooperation in Block 2 after a second-mover's history of play from Block 1 has exposed him as rational. The first two predictions are also consistent with the reputation-building theory of Kreps et al (1982), but the third is not. These results demonstrate that while the behavior we observe is inconsistent with Kreps et al (1982), it is rational in a reasonable non-equilibrium sense. While this is not the only model that could explain our experimental data, the recent evidence from Kagel and McGee (2014) indicates that this interpretation is worthy of a more careful exploration.

For a formal description of the model, let rounds be counted backwards. Play begins in round 10 and ends after round 1. We assume that all players adopt a strategy from  $S = \{s_{11}, s_{10}, \dots, s_1\}$ . A player adopting strategy  $s_k$  plays tit-for-tat in rounds 10 through  $k$  and defects in rounds  $k - 1$  through 1. Strategy  $s_{11}$  is defined as defecting in every round. In addition to the Kagel and McGee (2014) chat evidence, Embrey, Fréchette, and Yuksel (2014) find that subjects learn with experience to play exactly these strategies, so we believe that restricting the strategy space in this way has empirical foundations.

Players have prior beliefs  $\mu$  over their opponent's strategies in  $S$ . Thus,  $\mu(s_k)$  is the probability of playing against an opponent using strategy  $s_k$ . Though  $s_1$  is dominated for a payoff-maximizing

<sup>29</sup>This approach is similar to Radner (1986), in which players have arbitrary beliefs about the opponent's trigger strategy choice in a simultaneous-move FRPD and choose a best-response trigger strategy given these beliefs.

<sup>30</sup>Selten and Stoecker (1986) propose an alternative non-Bayesian model of learning from histories of play in FRPDs which predicts a general pattern of behavior that is consistent with our data. In their model, a player defects one period earlier (or later) with some probability if her previous opponent defected earlier (or later) than she did, and she defects in the same period otherwise. This learning model does not include beliefs about other players nor does it assume optimizing behavior, but only an iterative Markov-transition learning rule given a starting point and supergame outcome. Unlike our experiment, subjects in Selten and Stoecker (1986) are given no information about opponents' histories of play in prior supergames, and it is not clear how strategies would be updated in their model when players see the current opponent's history of play against others. In contrast to Selten and Stoecker, we model players as Bayesian optimizers in a framework that is general enough to accommodate the informational environment of our experiment as well as most other FRPD experiments.

<sup>31</sup>Evidence of non-equilibrium behavior like this is abundant in the experimental literature on strategic sophistication. See Crawford et al (2013) for a recent survey.

<sup>32</sup>This observation implies that cooperation would be sustained at least as long for more optimistic beliefs.

agent, we find a non-negligible amount of last-round cooperation in our data. Thus, we assume beliefs  $\mu$  are such that  $\mu(s_k) \in (0, 1)$  for all  $k$ , including  $k = 1$ .<sup>33</sup>

Players' beliefs are updated within each supgame round-by-round according to Bayes' rule based on the prior  $\mu$  and the opponents' history of actions in that supgame. If her opponent has cooperated up to and including round  $t + 1$ , a player believes that her opponent will continue playing tit-for-tat for at least one more round with probability  $p_t = \frac{\sum_{i=1}^t \mu(s_i)}{\sum_{i=1}^{t+1} \mu(s_i)}$ . In Proposition 1, we characterize the beliefs for which a naïve player chooses strategy  $s_k$  in terms of the conditional probability  $p_t$  that the opponent will continue to play tit-for-tat in round  $t$ , given that he has played tit-for-tat for all previous rounds, for all  $t$  up to round  $k$ .

**Proposition 1** (a) *The first-mover plays  $s_k$  if and only if*

$$p_l \geq \frac{4}{\sum_{i=k+1}^l (3 \prod_{j=i}^{l-1} p_j) + 7 \prod_{i=k}^{l-1} p_i}$$

for every  $l \in \{k, \dots, 10\}$ .

(b) *The second-mover plays  $s_{k+1}$  if and only if*

$$p_l \geq \max \left\{ \frac{1}{3}, \frac{5}{\sum_{i=k+1}^l (3 \prod_{j=i}^{l-1} p_j) + 8 \prod_{i=k}^{l-1} p_i} \right\}$$

for all  $l \in \{k, \dots, 10\}$ .

Proposition 1 provides lower bounds on the beliefs needed to sustain a particular strategy,  $s_k$ . For each round  $l \geq k$ , the subjective probability that the second-mover will play tit-for-tat until round  $l$  must be high enough that the expected payoff of cooperating in round  $l$  exceeds the payoff that can be obtained by defecting in round  $l$ . To build intuition, consider the condition for round  $l = k + 1$  given strategy  $s_k$ . First notice that the first-mover can always defect in rounds  $k + 1$  and  $k$  and obtain a total payoff of 8 from these two rounds (as a rational second-mover would respond by defecting, earning each player a payoff of 4 in each round). By playing tit-for-tat through round  $k$  instead, the first-mover faces three possible outcomes assuming that the second-mover has also chosen a strategy in  $S$ . She may earn payoffs of 0 in round  $k + 1$  and 4 in round  $k$  (if the second-mover defects in rounds  $k + 1$  and  $k$ ), payoffs of 7 in round  $k + 1$  and 0 in round  $k$  (if the second-mover cooperates in round  $k + 1$  and defects in round  $k$ ), or payoffs of 7 in both rounds  $k + 1$  and  $k$  (if the second-mover cooperates in both rounds). Thus, the first-mover is guaranteed a payoff of at least 4 from these two rounds. She can gain an additional 4 with certainty by defecting in rounds  $k + 1$  and  $k$ , but expects that she can gain either an additional 3 or 10 with some probability by playing tit-for-tat in rounds  $k + 1$  and  $k$ . If the subjective probability of these cooperation payoffs is sufficiently high, then it is rational for the first-mover to play tit-for-tat in round  $k + 1$ .

The second-mover's strategy is governed by similar belief conditions to the first-mover's. Consider the second-mover's condition for round  $l = k + 1$  given strategy  $s_k$  (assuming that the first-mover cooperates in round  $k + 1$ ). The second-mover can always defect in rounds  $k + 1$ ,  $k$  and  $k - 1$  to obtain a total payoff of 20 from these three rounds (12 in round  $k + 1$  and 4 in each of the following two rounds). By playing tit-for-tat through round  $k$  instead, the second-mover faces three possible

<sup>33</sup>We restrict  $\mu(s_k) \in (0, 1)$  so that Bayes' rule can always be used. Without changing the results of this section, we could instead assume players update via Bayes' rule whenever possible and allow beliefs to be free when zero-probability events are observed and an opponent's history does not eliminate any strategies and assign a new belief of zero when a particular strategy may be eliminated based on the opponent's history.

outcomes assuming that the first-mover has also chosen a strategy in  $S$ . She may earn payoffs of 7 in round  $k + 1$  and 4 in rounds  $k$  and  $k - 1$  (if the first-mover defects in rounds  $k$  and  $k - 1$ ), payoffs of 7 in rounds  $k + 1$  and  $k$  and 4 in round  $k - 1$  (if the first-mover cooperates in round  $k$  and defects in round  $k - 1$ ), or payoffs of 7 in rounds  $k + 1$  and  $k$  and 12 in round  $k - 1$  (if the first-mover cooperates in rounds  $k$  and  $k - 1$ ). Thus, the second-mover is guaranteed a payoff of at least 15 from these three rounds. She can gain an additional 5 with certainty by defecting in rounds  $k + 1$ ,  $k$  and  $k - 1$ , but expects that she can gain either an additional 3 or 11 with some probability by playing tit-for-tat in rounds  $k + 1$  and  $k$ . If the subjective probability of these cooperation payoffs is sufficiently high, then it is rational for the second-mover to play tit-for-tat in round  $k + 1$ . However, if the subjective probability that the first-mover plays tit-for-tat in any round is less than  $\frac{1}{3}$ , it is strictly dominant for the second-mover to defect before that round due to the incentive to defect before the first-mover (and earn a payoff of 12 instead of 4 for one round).

The conditions in Proposition 1 are permissive enough that cooperation is sustained into later rounds with a large variety of beliefs. The following examples demonstrate the range of beliefs that can support late-round cooperation.

*Example 1 (Uniform Prior)* Assume that both players have prior beliefs such that  $\mu(s_k) = 1/11$  for all  $k$ . Then by Proposition 1 the first-mover's optimal Block 1 strategy is  $s_2$  and the second-mover's optimal Block 1 strategy is  $s_3$ .<sup>34</sup> That is, the first-mover plays tit-for-tat until the last round, in which she always defects, while the second-mover plays tit-for-tat up to the next-to-last round and always defects in the last two rounds. Hence, cooperation until the penultimate round is observed.

*Example 2 (Pessimistic Triangular Prior)* Assume that both players have prior beliefs  $\mu(s_k) = k/66$  for all  $k$ . Then by Proposition 1 the first-mover's optimal Block 1 strategy is  $s_3$  and the second-mover's optimal Block 1 strategy is  $s_5$ . That is, the first-mover plays tit-for-tat for eight rounds and then defects, while the second-mover plays tit-for-tat for six rounds and then defects thereafter. These relatively pessimistic beliefs still support cooperation for more than half of the supergame.

Now consider how players update their beliefs based on revealed Block 1 histories. We assume that players' own past opponents' behavior from either Blocks 1 or 2 does not affect beliefs because players know that they will not face previous opponents again. In Block 2, however, players incorporate their opponents' histories into their beliefs when available. Let  $\tilde{\mu}$  represent the updated beliefs based on the prior  $\mu$  and the observed Block 1 history. Because players are assumed to adopt a pure strategy from  $S$ , beliefs are updated such that if a player's opponent never defected before her opponent in rounds  $10, \dots, n$ , then  $\sum_{i=n-1}^{11} \tilde{\mu}(s_i) = 0$  holds. If the opponent always defected before her opponent by round  $m$ , then  $\sum_{i=1}^m \tilde{\mu}(s_i) = 0$  holds. If the opponent's opponent always defected first, then  $\tilde{\mu}(s_i) = \mu(s_i)$  holds for all  $i$ .

If players focus on how long they can expect their opponent to cooperate, then first-movers could have a fairly optimistic prior (i.e., first-movers may be trusting in Block 1) and then, upon learning that their opponent was a Defector in Block 1, become more pessimistic and cooperate less in Block 2. This argument is formalized in the following proposition.

<sup>34</sup>Calculation of the conditional probabilities  $p_t$  and the conditions in Proposition 1 for a uniform prior show that this is the case. Take the first-mover for example. Because  $p_t$  decreases in  $t$  for these beliefs,  $s_k$  is optimal if and only if  $p_k \geq 4/7$  holds because  $p_k \geq 4/7$  implies that  $p_l \geq 4 / \left[ \sum_{j=k+1}^l (3 \prod_{j=i}^{l-1} p_j) + 7 \prod_{i=k}^{l-1} p_i \right]$  holds for all  $l \geq k$ .  $s_2$  is optimal because  $p_k \geq 4/7$  holds for all  $k \geq 2$  and  $p_1 < 4/7$ . The calculation is similar for the second-mover.



**Proposition 2** *Suppose that the second-mover always defected before her opponent by round  $n$  of Block 1 supergames. If the naïve first-mover's prior beliefs satisfy  $\mu(s_{k+1}) \leq (3/4) \sum_{i=1}^k \mu(s_i)$  for all  $k \in \{m, \dots, 10\}$ , where  $m \leq n$ , then her Block 1 strategy is  $s_m$  and her Block 2 strategy is  $s_{m+t}$  for some  $t \geq 1$ .*

Proposition 2 applies to the set of first-mover prior beliefs such that, for a certain number of rounds beginning with the first, the probability that the second-mover will begin unconditionally defecting in each of these rounds is not more than three-fourths the probability that the second-mover will play tit-for-tat in that round. This set includes Examples 1 and 2 above and infinitely many others. Given such beliefs, a first-mover who is classified as Trusting in Block 1, playing strategy  $s_i$ , will respond to a second-mover's Defector history ( $s_j, j \geq i$ ) by choosing a strategy  $s_{i+t}, t \geq 1$  in Block 2. In other words, if the first-mover had been playing tit-for-tat up to round  $i$  in Block 1, she will defect at least one round earlier in Block 2 games against a second-mover whose game history shows that he always defected in round  $i$  or earlier in Block 1. Again, this prediction is consistent with our data, which shows a significant decrease in cooperation by Trusting first-movers whose opponents are revealed as Defectors compared to those whose opponents' types are not revealed. Proposition 2 also predicts earlier defections in a given supergame by Trusting first-movers when the second-mover is revealed to be a Defector than when no information is revealed, as observed in the data.

In contrast, a first-mover may become more optimistic about how long he can expect to cooperate when his opponent is revealed to be an Imitator through her Block 1 history. In this way, a first-mover will choose to cooperate longer upon seeing that his opponent was cooperative in Block 1. The following proposition formalizes this argument.

**Proposition 3** *Suppose that the second-mover never defected before her opponent in rounds  $10, \dots, m$  of Block 1 supergames.*

- (a) *If the naïve first-mover's prior beliefs satisfy  $\mu(s_{k+1}) \leq (3/4) \sum_{i=1}^k \mu(s_i)$  for all  $k \in \{n, \dots, 10\}$ , where  $n > m$ , then her Block 1 strategy is  $s_n$  and her Block 2 strategy is  $s_{n-t}$  for some  $t \geq 1$ .*
- (b) *If the naïve first-mover's prior beliefs satisfy  $\mu(s_{11}) > 3/7$  then her Block 1 strategy is  $s_{11}$  and her Block 2 strategy is  $s_{11-t}$  for some  $t \geq 1$ .*

Proposition 3 applies to the same set of prior beliefs as Proposition 2 as well as beliefs under which the first-mover defects in every round of Block 1 supergames. Given these beliefs, a first-mover plays strategy  $s_{11}$ , is classified as Non-Trusting in Block 1, and responds to a second-mover's Imitator-type history ( $s_i, i \leq 10$ ) by choosing a strategy  $s_j, j \leq 10$  in Block 2. The simple intuition for this result is that if the first-mover had been playing tit-for-tat up to round  $i$  in Block 1, she will continue to play tit-for-tat at least one round later in Block 2 games against a second-mover whose game history shows that he always played tit-for-tat beyond round  $i$  in Block 1. This prediction is consistent with our data, which shows a significant increase in initial cooperation by Non-Trusting first-movers when their opponents are revealed as Imitators compared to those whose opponents' types are not revealed, a finding that is clearly inconsistent with the predictions of the Kreps et al (1982) model. This model may also provide a more plausible explanation than Kreps et al (1982) for the finding of Gächter and Thoni (2005)'s public goods experiment, where low contributors contribute more when grouped with other players revealed to have made low contributions in the past.

## 7 CONCLUSION

We have shown that cooperation in FRPDs occurs when the reputation-building theory of Kreps et al (1982) predicts complete unraveling. The results of the experiment indicate that first-movers change their strategy when they observe their opponent's history of play by either increasing or decreasing their degree of cooperation based on the relative cooperativeness of their opponent. First-movers tend to cooperate at least as often initially and continue cooperating at least as long when second-mover histories are revealed, except in the case of relatively trusting first-movers meeting relatively uncooperative second-movers. Second-movers also tend to behave more cooperatively when their histories are revealed. In particular, we find the surprising result that revealing histories improves cooperation even in the case of a relatively Non-Trusting first-mover meeting a relatively uncooperative second-mover. Thus, cooperation persists and often increases, even when revealed histories are relatively uncooperative. These results are clearly inconsistent with the reputation-building theory.

We show that an alternative behavioral model to Kreps et al (1982) generates predictions that are consistent with the features of our experimental data. Players in this model form beliefs over the strategies of their opponents which may not be consistent with the opponent's best response, and then choose the optimal strategy based on those naïve beliefs. We do not view this as the ultimate model of behavior in FRPDs, but as a simple and reasonable one which generates predictions that fit the observed behavior better than prevailing equilibrium models. By using such a simple model, we avoid *ad hoc* assumptions about more specific behavioral types which could possibly fit behavior in this game more precisely. One limitation of this analysis is that beliefs are a critical part of our behavioral model, but we are able to observe beliefs only in a very limited way. Because our main hypotheses could be tested without elicited beliefs, and because eliciting beliefs before or during gameplay is complicated and may itself alter beliefs and behavior, we opted not to do so. Examining beliefs in more depth may be an interesting direction for future research.

**Acknowledgements** The authors thank Yaron Azrieli, Lucas Coffman, Glenn Dutcher, John Kagel, Semin Kim, Peter McGee, Xiangyu Qu, Arno Riedl, Dan Schley, Mike Sinkey, Tom Wilkening, and Chao Yang for their helpful comments and conversations. Healy acknowledges financial support from National Science Foundation grant #SES-0847406. Any opinions, findings, conclusions or recommendations expressed are those of the authors and do not necessarily reflect the views of the Federal Trade Commission.

## APPENDIX A

**Table A.1** First and second-movers' Frequency of Defecting First

		Trusting first-mover					
		2S			1S		
		1st	2nd		1st	2nd	
2nd Mover		Mover	Mover	Neither	Mover	Mover	Neither
Imitator		68.4%	31.6%	0.0%	42.9%	57.1%	0.0%
Cooperator		43.6%	40.4%	16.0%	53.2%	29.9%	16.9%
Defector		63.0%	37.0%	0.0%	38.5%	50.8%	10.8%

		Non-Trusting first-mover					
		2S			1S		
		1st	2nd		1st	2nd	
2nd Mover		Mover	Mover	Neither	Mover	Mover	Neither
Imitator		81.8%	18.2%	0.0%	84.4%	12.5%	3.1%
Cooperator		62.7%	21.6%	15.7%	77.8%	12.7%	9.5%
Defector		63.2%	31.6%	5.3%	76.7%	18.3%	5.0%

**Table A.2** Mean number of round before first defection (conditional defection) for first-movers (second-movers)

		Trusting first-mover			
		2S		1S	
2nd Mover		1st Mover	2nd Mover	1st Mover	2nd Mover
Imitator		7.6	8.1	8.2	8.3
Cooperator		8.1	8.2	6.6	8.1
Defector		3.2	4.1	6.4	6.3

		Non-Trusting first-mover			
		2S		1S	
2nd Mover		1st Mover	2nd Mover	1st Mover	2nd Mover
Imitator		4.0	7.5	2.2	5.0
Cooperator		4.9	7.6	2.9	6.3
Defector		3.5	4.5	1.2	3.4

**Table A.3** Elicited Beliefs

	Predicted # of Rounds of Cooperation	
	1S	2S
<b>Respondent is a first-mover</b>		
× first-mover avg. rounds coop. in block 1	0.497 ( $< 0.001$ )	-0.011 (0.920)
× second-mover avg. rounds coop. in block 1		0.113 (0.769)
<b>Respondent is a second-mover</b>		
× first-mover avg. rounds coop. in block 1	0.447 (0.665)	-7.134 (0.098)
× second-mover avg. rounds coop. in block 1	0.129 (0.339)	0.092 (0.642)
<b>Constant</b>	3.201 ( $< 0.001$ )	6.689 (0.042)
<b>Observations</b>	460	200

*Notes:* p-values from Robust Standard Errors Clustered by Individual Mover.

**Table A.4** The Effect of Player Types on First-Round Cooperation and Rounds of Cooperation before First Defection – Types Defined by Last Block 1 Supergame

	1st Round Cooperation by First-Mover		# of Rounds of Cooperation Before 1st Defection			
	(i)	(ii)	(iii)	(iv)	(v)	(vi)
Avg. Block 1 Coop.	0.042 (0.251)		0.455 (0.006)			
× S2	0.116 (0.013)		0.253 (0.169)			
Defector		-0.274 (0.489)				
Imitator		0.706 (0.113)		1.264 (0.377)	2.698 (< 0.001)	0.996 (0.588)
Cooperator		1.127 (0.010)		2.278 (0.006)	1.950 (0.038)	2.425 (0.016)
Trusting						
× Defector		1.398 (0.005)		8.374 (< 0.001)		
× Imitator		2.796 (< 0.001)		9.921 (< 0.001)		
× Cooperator		2.114 (< 0.001)		8.144 (< 0.001)		
S2						
× Defector				4.701 (< 0.001)	-4.950 (< 0.001)	5.729 (< 0.001)
× Imitator				4.434 (0.061)	0.040 (0.939)	6.079 (0.042)
× Cooperator				5.028 (0.004)	0.953 (0.556)	6.521 (0.004)
S2×Trusting						
× Defector				-10.205 (< 0.001)		
× Imitator				-4.520 (0.039)		
× Cooperator				-4.042 (0.003)		
Supergame#	0.385 (0.035)	0.493 (0.041)	0.760 (0.129)	0.839 (0.143)	-0.402 (0.466)	4.517 (0.019)
(Supergame#) <sup>2</sup>	-0.047 (0.121)	-0.061 (0.126)	-0.081 (0.300)	-0.099 (0.262)	0.066 (0.469)	-0.596 (0.049)
Constant	-0.307 (0.226)	-1.223 (0.001)	-0.321 (0.653)	-4.219 (< 0.001)	5.983 (< 0.001)	-11.873 (< 0.001)
Observations	660	660	660	660	405	255

Notes: Column (v) includes observations with Trusting first-mover and column (vi) includes observations with Non-Trusting first-mover only. p-values from Robust Standard Errors Clustered by Individual 1st Mover.

**Proposition 4** Let  $p \in (0, 1)$  be the common belief that the other player plays tit-for-tat and  $\bar{p}_t$  the period  $t$  posterior belief. The following is a sequential equilibrium for a sequential-move FRPD.

(a) The second-mover plays tit-for-tat in round  $t$  with probability

$$q_t(p) = \min \left\{ \frac{p}{1-p} \frac{1-\bar{p}_{t+1}}{\bar{p}_{t+1}}, 1 \right\}.$$

Otherwise, the second-mover defects in round  $t$ .

(b) The first-mover cooperates in round  $t$  if and only if  $t \geq t^*(p)$  where

$$t^*(p) = \min \{t \in \mathbb{N} : p \geq \bar{p}_t\} \text{ and } \bar{p}_t = \left(\frac{4}{7}\right)^t \text{ hold for all } t.$$

Otherwise, the first-mover defects in round  $t$ .

*Proof* It is easier notationally to derive the equilibrium by counting backwards with  $t = 10$  representing the first round of the supreme. In the body of the paper, however, time is indexed forward with  $t = 1$  representing the first round of the supreme. Now, in any period  $t$ , the first-mover will cooperate if

$$\begin{aligned} p_t \left( 7 + V_{t-1} \left( \frac{p_t}{p_t + (1-p_t)q_t(p_t)} \right) \right) \\ + (1-p_t) \left[ q_t(p_t) \left( 7 + V_{t-1} \left( \frac{p_t}{p_t + (1-p_t)q_t(p_t)} \right) \right) + (1-q_t(p_t))V_{t-1}(0) \right] \\ \geq 4 + V_{t-1}(p_t), \end{aligned}$$

where  $V_{t-1}(p)$  is the continuation value of the first-mover entering period  $t-1$  with belief  $p$ . Let  $V_0 \equiv 0$ . Let  $\bar{p}_t$  be the smallest value of  $p_t$  satisfying this inequality. (We will show later that the inequality in fact grows in  $p_t$ .)

The probability a selfish second-mover cooperates is the highest  $q$  such that the first-mover is willing to cooperate in periods  $t-1$  after observing cooperation in period  $t$ . Thus, if  $\bar{p}_t$  is the lowest belief at which first-mover will cooperate in period  $t$ , then  $q_t(p)$  solves

$$q_t(p) = \arg \max \left\{ q \in [0, 1] : \frac{p}{p+q(1-p)} \geq \bar{p}_{t-1} \right\},$$

and so

$$q_t(p) = \min \left\{ \frac{p}{1-p} \frac{1-\bar{p}_{t-1}}{\bar{p}_{t-1}}, 1 \right\}.$$

For completeness, let  $q_t(1) = 1$  for all  $t$  and  $q_t \equiv 1$  for any  $t$  where  $\bar{p}_{t-1} = 0$ . Since a selfish second-mover never cooperates in the last period, set  $q_1(p) = 0$  for all  $p$ . (This is equivalent to setting  $\bar{p}_0 = 1$ .)

For any  $t > 1$ , consider the case where  $p_t \geq \bar{p}_{t-1}$ . Here,  $q_t(p_t) = 1$  (the second-mover cooperates with certainty) and

$$\frac{p_t}{p_t + (1-p_t)q_t(p_t)} = p_t,$$

‘so the above inequality becomes

$$p_t(7 + V_{t-1}(p_t)) + (1-p_t)[(7 + V_{t-1}(p_t))] \geq 4 + V_{t-1}(p_t),$$

or

$$7 \geq 4.$$

In other words, the first-mover always cooperates if  $p_t \geq \bar{p}_{t-1}$ . This proves that  $\bar{p}_t \leq \bar{p}_{t-1}$ .

Now suppose  $p_t < \bar{p}_{t-1}$ . Here

$$q_t(p_t) = \frac{p_t}{1-p_t} \frac{1-\bar{p}_{t-1}}{\bar{p}_{t-1}}$$

and so

$$\frac{p_t}{p_t + (1-p_t)q_t(p_t)} = \bar{p}_{t-1}.$$

The above inequality becomes

$$\begin{aligned} + (1-p_t) \left[ \frac{p_t}{1-p_t} \frac{1-\bar{p}_{t-1}}{\bar{p}_{t-1}} (7 + V_{t-1}(\bar{p}_{t-1})) + \left( 1 - \frac{p_t}{1-p_t} \frac{1-\bar{p}_{t-1}}{\bar{p}_{t-1}} \right) V_{t-1}(0) \right] \\ \geq 4 + V_{t-1}(p_t). \end{aligned}$$

After several steps of algebra, this reduces to

$$p_t \geq \bar{p}_{t-1} \frac{4 + V_{t-1}(p_t) - V_{t-1}(0)}{7 + V_{t-1}(\bar{p}_{t-1}) - V_{t-1}(0)}.$$

Since  $\bar{p}_t$  solves this inequality exactly, it has the property that

$$\bar{p}_t = \bar{p}_{t-1} \frac{4 + V_{t-1}(\bar{p}_t) - V_{t-1}(0)}{7 + V_{t-1}(\bar{p}_{t-1}) - V_{t-1}(0)}.$$

In  $t = 1$  the first-mover cooperates if

$$p_1(7) + (1 - p_1)[0] \geq 4,$$

and so

$$\bar{p}_1 = 4/7.$$

Thus,

$$V_1(p) = \begin{cases} 7p & \text{if } p \geq \bar{p}_1 \\ 4 & \text{otherwise} \end{cases},$$

or

$$V_1(p) = \max\{7p, 4\}.$$

Note that  $V_1(p_1) = 4 = V_1(0)$  for any  $p_1 \leq \bar{p}_1$ .

In  $t = 2$  we know that if  $p_2 \geq \bar{p}_1 = 4/7$  then the first-mover cooperates with certainty.

If  $p_2 < \bar{p}_1$  then he will cooperate only if  $p_2 \geq \bar{p}_2$ , where  $\bar{p}_2$  solves

$$\begin{aligned} \bar{p}_2 &= \bar{p}_1 \frac{4 + V_1(\bar{p}_2) - V_1(0)}{7 + V_1(\bar{p}_1) - V_1(0)} \\ &= \left(\frac{4}{7}\right)^2. \end{aligned}$$

The expression for  $V_2(p)$  is given by

$$V_2(p) = \begin{cases} p(7 + V_1(p)) + (1 - p)(7 + V_1(p)) & \text{if } p \geq \bar{p}_1 \\ p(7 + V_1(\bar{p}_1)) + (1 - p) \left( \frac{p}{1-p} \frac{1-\bar{p}_1}{\bar{p}_1} (7 + V_1(\bar{p}_1)) + \left(1 - \frac{p}{1-p} \frac{1-\bar{p}_1}{\bar{p}_1}\right) (V_1(0)) \right) & \text{if } p \in [\bar{p}_2, \bar{p}_1) \\ 4 + V_1(p) & \text{otherwise} \end{cases},$$

which is equal to

$$V_2(p) = \begin{cases} 7 + 7p & \text{if } p \geq \bar{p}_1 \\ 7 \frac{p}{\bar{p}_1} + 4 & \text{if } p \in [\bar{p}_2, \bar{p}_1) \\ 4 + 4 & \text{otherwise} \end{cases}.$$

Note that  $V_2(p_2) = 8 = V_2(0)$  for any  $p_2 \leq \bar{p}_2$ .

In  $t = 3$  we know that if  $p_3 \geq \bar{p}_2$  then the first-mover cooperates with certainty.

If  $p_3 < \bar{p}_2$  then he will cooperate only if  $p_3 \geq \bar{p}_3$ , where  $\bar{p}_3$  solves

$$\begin{aligned} \bar{p}_3 &= \bar{p}_2 \frac{4 + V_2(\bar{p}_3) - V_2(0)}{7 + V_2(\bar{p}_2) - V_2(0)} \\ &= \left(\frac{4}{7}\right)^3. \end{aligned}$$

The expression for  $V_3(p)$  is given by

$$V_3(p) = \begin{cases} p(7 + V_2(p)) + (1 - p)(7 + V_2(p)) & \text{if } p \geq \bar{p}_2 \\ p(7 + V_2(\bar{p}_2)) + (1 - p) \left( \frac{p}{1-p} \frac{1-\bar{p}_2}{\bar{p}_2} (7 + V_2(\bar{p}_2)) + \left(1 - \frac{p}{1-p} \frac{1-\bar{p}_2}{\bar{p}_2}\right) (V_2(0)) \right) & \text{if } p \in [\bar{p}_3, \bar{p}_2) \\ 4 + V_2(p) & \text{otherwise} \end{cases},$$

which is equal to

$$V_3(p) = \begin{cases} 7 + 7 + 7p & \text{if } p \geq \bar{p}_2 \\ 7 + 7 \frac{p}{\bar{p}_2} + 4 & \text{if } p \in [\bar{p}_3, \bar{p}_2) \\ 7 \frac{p}{\bar{p}_2} + 4 + 4 & \text{if } p \in [\bar{p}_3, \bar{p}_2) \\ 4 + 4 + 4 & \text{otherwise} \end{cases}.$$

In general we will have

$$\bar{p}_t = \left(\frac{4}{7}\right)^t$$

and

$$q_t(p) = \min\left\{\frac{p}{1-p} \frac{1-\bar{p}_{t-1}}{\bar{p}_{t-1}}, 1\right\}.$$

Let

$$t^*(p) = \min\{t \in \mathbb{N} : p \geq \bar{p}_t\}.$$

First-movers will cooperate in period  $T$  if  $p_T^* \geq \bar{p}_T$ . Thereafter they will cooperate as long as they have never seen a defection and will never cooperate after seeing a defection. In that case beliefs will evolve according to the formula

$$p_t^* = \begin{cases} p_{t+1}^* & \text{if } t \geq t^*(p_T^*) \\ \bar{p}_t & \text{otherwise} \end{cases}.$$

Beliefs change to  $p_t^* = 0$  if a defection is observed in any previous period. If  $p_T^* < \bar{p}_T$  then both players always defect and  $p_t^* = p_T^*$  for every period  $t$ . The on-path continuation value of the first-mover will equal

$$V_t(p) = \begin{cases} 7(t - t^*(p)) + 7\frac{p}{\bar{p}_{t^*(p)-1}} + 4(t^*(p) - 1) & \text{if } t^*(p) \leq t \\ 4t & \text{if } t^*(p) > t \end{cases},$$

where we set  $\bar{p}_0 = 1$ .

#### Proof of Proposition 1

(a) Let the first-mover's expected payoff in round  $s$  from the remaining rounds  $1, \dots, s$  given beliefs  $p_1, \dots, p_s$  be denoted by  $V_s(p_1, \dots, p_s)$ . The expected payoff for cooperating in round  $t$  is  $p_t(7 + V_{t-1}(p_1, \dots, p_{t-1})) + (1 - p_t)V_{t-1}(0, \dots, 0)$ . The expected payoff for defecting in round  $t$ , given that the second-mover will respond by defecting for at least one round, is at most  $4 + V_{t-1}(p_1, \dots, p_t p_{t-1})$ . Therefore, the first-mover plays tit-for-tat in period  $t$  if and only if the following inequality holds:

$$p_t(7 + V_{t-1}(p_1, \dots, p_{t-1})) + (1 - p_t)V_{t-1}(0, \dots, 0) \geq 4 + V_{t-1}(p_1, \dots, p_t p_{t-1}).$$

We need to show that:

$$V_t(p_1, p_2, \dots, p_t) = 4(t-1) + \sum_{i=k+1}^t (3 \prod_{j=i}^t p_j) + 7 \prod_{i=k}^t p_i$$

if

$$p_l \geq \frac{4}{\sum_{i=k+1}^l (3 \prod_{j=i}^{l-1} p_j) + 7 \prod_{i=k}^{l-1} p_i} \quad \forall t \geq l \geq k$$

and

$$V_t(p_1, p_2, \dots, p_t) = 4t$$

otherwise.

The proof is by induction. First, we know that the first-mover cooperates in round 1 if and only if  $p_1 7 + (1 - p_1)0 \geq 4$  holds. Therefore, if  $p_1 \geq \frac{4}{7}$  holds then we have  $V_1(p_1) = 7p_1$ , and if  $p_1 < \frac{4}{7}$  holds then we have  $V_1(p_1) = 4$ , and the formula is true for  $t = 1$ .

Now, assume that the formula holds for all rounds up to  $t-1$  and show that it holds for round  $t$ . Assume that the following holds:

$$V_{t-1}(p_1, p_2, \dots, p_{t-1}) = 4(t-2) + \sum_{i=k+1}^{t-1} (3 \prod_{j=i}^{t-1} p_j) + 7 \prod_{i=k}^{t-1} p_i$$

if

$$p_l \geq \frac{4}{\sum_{i=k+1}^l (3 \prod_{j=i}^{l-1} p_j) + 7 \prod_{i=k}^{l-1} p_i} \quad \forall t-1 \geq l \geq k$$

and

$$V_{t-1}(p_1, p_2, \dots, p_{t-1}) = 4(t-1)$$

otherwise.

The first-mover cooperates in round  $t$  if and only if:

$$p_t(7 + V_{t-1}(p_1, \dots, p_{t-1})) + (1 - p_t)V_{t-1}(0, \dots, 0) \geq 4 + V_{t-1}(p_1, \dots, p_t p_{t-1}).$$



We have assumed that  $p_l \geq \frac{4}{\sum_{i=k+1}^l (3 \prod_{j=i}^{l-1} p_j) + 7 \prod_{i=k}^{l-1} p_i}$  holds for all  $l$  such that  $t-1 \geq l \geq k$ . First, suppose also that  $p_t p_{t-1} \geq \frac{4}{\sum_{i=k+1}^{t-1} (3 \prod_{j=i}^{t-2} p_j) + 7 \prod_{i=k}^{t-2} p_i}$  holds. Then the first-mover cooperates in round  $t$  if and only if:

$$\begin{aligned} & p_t(7+4(t-2)) + \sum_{i=k+1}^{t-1} (3 \prod_{j=i}^{t-1} p_j) + 7 \prod_{i=k}^{t-1} p_i + (1-p_t)4(t-1) \geq 4 + 4(t-2) + \sum_{i=k+1}^{t-1} (3 \prod_{j=i}^t p_j) + 7 \prod_{i=k}^t p_i \\ \Leftrightarrow & 7p_t + 4(t-2)p_t + \sum_{i=k+1}^{t-1} (3 \prod_{j=i}^t p_j) + 7 \prod_{i=k}^t p_i + 4(t-1) - 4(t-1)p_t \geq 4(t-1) + \sum_{i=k+1}^{t-1} (3 \prod_{j=i}^t p_j) + 7 \prod_{i=k}^t p_i \\ \Leftrightarrow & (7-4)p_t \geq 0 \Leftrightarrow p_t \geq 0. \end{aligned}$$

Now suppose that  $p_t p_{t-1} < \frac{4}{\sum_{i=k+1}^{t-1} (3 \prod_{j=i}^{t-2} p_j) + 7 \prod_{i=k}^{t-2} p_i}$  holds. Then the first-mover cooperates in round  $t$  if and only if:

$$\begin{aligned} & p_t(7+4(t-2)) + \sum_{i=k+1}^{t-1} (3 \prod_{j=i}^{t-1} p_j) + 7 \prod_{i=k}^{t-1} p_i + (1-p_t)4(t-1) \geq 4t \\ \Leftrightarrow & p_t(3 + \sum_{i=k+1}^{t-1} (3 \prod_{j=i}^{t-1} p_j) + 7 \prod_{i=k}^{t-1} p_i) \geq 4 \\ \Leftrightarrow & p_t \geq \frac{4}{\sum_{i=k+1}^t (3 \prod_{j=i}^{t-1} p_j) + 7 \prod_{i=k}^{t-1} p_i}. \end{aligned}$$

Hence, the first-mover cooperates in round  $t$  and  $V_t(p_1, p_2, \dots, p_t) = 4(t-1) + \sum_{i=k+1}^t (3 \prod_{j=i}^t p_j) + 7 \prod_{i=k}^t p_i$  if and only if  $p_l \geq \frac{4}{\sum_{i=k+1}^l (3 \prod_{j=i}^{l-1} p_j) + 7 \prod_{i=k}^{l-1} p_i}$  holds for all  $l$  such that  $t \geq l \geq k$ . Otherwise, the first-mover defects in round  $t$  and  $V_t(p_1, p_2, \dots, p_t) = 4t$ .

(b) Let the second-mover's expected payoff in round  $s$  from the remaining rounds  $1, \dots, s-1$  given beliefs  $p_1, \dots, p_{s-1}$  be denoted by  $V_s(p_1, \dots, p_{s-1})$ . The expected payoff for cooperating in round  $t$  is  $7 + p_t(7 + V_t(p_1, \dots, p_{t-1})) + (1-p_t)(4 + V_t(0, \dots, 0))$ . The expected payoff for defecting in round  $t$ , given that the first-mover will respond by defecting for at least one round, is at most  $12 + V_t(p_1, \dots, p_{t-1})$ . Therefore, the second-mover plays tit-for-tat in period  $t+1$  if and only if the following inequality holds:

$$7 + p_t(7 + V_t(p_1, \dots, p_{t-1})) + (1-p_t)(4 + V_t(0, \dots, 0)) \geq 12 + V_t(p_1, \dots, p_{t-1}).$$

We need to show that:

$$V_{t+1}(p_1, p_2, \dots, p_t) = 4(t-1) + \sum_{i=k+1}^t (3 \prod_{j=i}^t p_j) + 8 \prod_{i=k}^t p_i$$

if

$$p_l \geq \max\left\{\frac{1}{3}, \frac{5}{\sum_{i=k+1}^l (3 \prod_{j=i}^{l-1} p_j) + 8 \prod_{i=k}^{l-1} p_i}\right\} \forall t \geq l \geq k$$

and

$$V_{t+1}(p_1, p_2, \dots, p_t) = 4(t-1)$$

otherwise.

The proof is by induction. First, we know that defection is the dominant action for the second-mover in round 1. The second-mover cooperates in round 2 if and only if  $7 + p_1 12 + (1-p_1)4 \geq 12 + 4$  holds. Therefore, if  $p_1 \geq \frac{5}{8}$  holds then we have  $V_2(p_1) = 4 + 8p_1$ , and if  $p_1 < \frac{5}{8}$  holds then we have  $V_2(p_1) = 4$ , and the formula is true for  $t = 1$ .

Now, we assume that the formula holds for all rounds up to  $t-1$  and show that it holds for round  $t$ . Assume that the following holds:

$$V_t(p_1, p_2, \dots, p_{t-1}) = 4(t-2) + \sum_{i=k+1}^{t-1} (3 \prod_{j=i}^{t-1} p_j) + 8 \prod_{i=k}^{t-1} p_i$$

if

$$p_l \geq \frac{5}{\sum_{i=k+1}^l (3 \prod_{j=i}^{l-1} p_j) + 8 \prod_{i=k}^{l-1} p_i} \forall t-1 \geq l \geq k$$

and

$$V_t(p_1, p_2, \dots, p_{t-1}) = 4(t-2)$$

otherwise.

The second-mover cooperates in round  $t+1$  if and only if:

$$7 + p_t(7 + V_t(p_1, \dots, p_{t-1})) + (1 - p_t)(4 + V_t(0, \dots, 0)) \geq 12 + V_t(p_1, \dots, p_t p_{t-1}).$$

We have assumed that  $p_l \geq \frac{5}{\sum_{i=k+1}^l (3 \prod_{j=i}^{l-1} p_j) + 8 \prod_{i=k}^{l-1} p_i}$  holds for all  $l$  such that  $t-1 \geq l \geq k$ . First, suppose that  $p_t p_{t-1} \geq \frac{5}{\sum_{i=k+1}^{t-1} (3 \prod_{j=i}^{t-2} p_j) + 8 \prod_{i=k}^{t-2} p_i}$  holds. Then the second-mover cooperates in round  $t+1$  if and only if:

$$\begin{aligned} 7 + p_t(7 + 4(t-2) + \sum_{i=k+1}^{t-1} (3 \prod_{j=i}^{t-1} p_j) + 8 \prod_{i=k}^{t-1} p_i) + (1 - p_t)4(t-1) &\geq 12 + 4(t-2) + \sum_{i=k+1}^{t-1} (3 \prod_{j=i}^t p_j) + 8 \prod_{i=k}^t p_i \\ \Leftrightarrow 11 + 4(t-2) + 3p_t + \sum_{i=k+1}^{t-1} (3 \prod_{j=i}^t p_j) + 8 \prod_{i=k}^t p_i &\geq 12 + 4(t-2) + \sum_{i=k+1}^{t-1} (3 \prod_{j=i}^t p_j) + 8 \prod_{i=k}^t p_i \\ \Leftrightarrow p_t &\geq \frac{1}{3}. \end{aligned}$$

Now suppose that  $p_t p_{t-1} < \frac{5}{\sum_{i=k+1}^{t-1} (3 \prod_{j=i}^{t-2} p_j) + 8 \prod_{i=k}^{t-2} p_i}$  holds. Then the second-mover cooperates in round  $t+1$  if and only if:

$$\begin{aligned} 7 + p_t(7 + 4(t-2) + \sum_{i=k+1}^{t-1} (3 \prod_{j=i}^{t-1} p_j) + 8 \prod_{i=k}^{t-1} p_i) + (1 - p_t)4(t-1) &\geq 12 + 4(t-1) \\ \Leftrightarrow p_t(3 + \sum_{i=k+1}^{t-1} (3 \prod_{j=i}^{t-1} p_j) + 8 \prod_{i=k}^{t-1} p_i) &\geq 5 \\ \Leftrightarrow p_t &\geq \frac{5}{\sum_{i=k+1}^t (3 \prod_{j=i}^{t-1} p_j) + 8 \prod_{i=k}^{t-1} p_i}. \end{aligned}$$

Hence, the second-mover cooperates in round  $t+1$  and  $V_{t+1}(p_1, p_2, \dots, p_t) = 4(t-1) + \sum_{i=k+1}^t (3 \prod_{j=i}^t p_j) + 8 \prod_{i=k}^t p_i$  if  $p_l \geq \max\{\frac{1}{3}, \frac{5}{\sum_{i=k+1}^l (3 \prod_{j=i}^{l-1} p_j) + 8 \prod_{i=k}^{l-1} p_i}\}$  holds for all  $l$  such that  $t \geq l \geq k$ . Otherwise, the second-mover defects in round  $t+1$  and  $V_{t+1}(p_1, p_2, \dots, p_t) = 4t$ .

#### Proof of Proposition 2

By Proposition 1, the first-mover's Block 1 strategy is  $s_m$  if and only if  $\mu$  is such that  $p_k \geq \frac{4}{\sum_{i=m+1}^k (3 \prod_{j=i}^{k-1} p_j) + 7 \prod_{i=m}^{k-1} p_i}$  holds for all  $k \in \{m, \dots, 10\}$ . This condition can be re-written in terms of the prior beliefs  $\mu$  as:

$$\frac{\sum_{i=1}^k \mu(s_i)}{\sum_{i=1}^{k+1} \mu(s_i)} \geq \frac{4}{\sum_{i=m+1}^k (3 \prod_{j=i}^{k-1} (\sum_{l=1}^j \mu(s_l)) / \sum_{l=1}^{j+1} \mu(s_l)) + 7 \prod_{i=m}^{k-1} (\sum_{l=1}^i \mu(s_l)) / \sum_{l=1}^{i+1} \mu(s_l)}$$

for all  $k \in \{m, \dots, 10\}$ . After several steps of algebra, the denominator of the right-hand-side of the above inequality simplifies to  $\frac{1}{\sum_{i=1}^k \mu(s_i)} ((7 + 3(k-n)) \sum_{i=1}^m \mu(s_i) + 3 \sum_{i=m+1}^k ((k+1-i) \mu(s_i)))$ , and the condition can be simplified to:

$$\mu(s_{k+1}) \leq \frac{3}{4}(k+1-m) \sum_{i=1}^m \mu(s_i) + \frac{1}{4} \sum_{i=m+1}^k ((3(k-i)-1) \mu(s_i))$$

for all  $k \in \{m, \dots, 10\}$ .

Now suppose that the first-mover's prior beliefs satisfy  $\mu(s_{k+1}) \leq (3/4) \sum_{i=1}^k \mu(s_i)$  for all  $k \in \{m, \dots, 10\}$ . For  $k = m$ , the above condition is satisfied trivially. We now show that the above condition is satisfied for  $k = m+r$  for any  $r \geq 1$ . For any  $r \geq 1$ , the inequality  $\mu(s_{m+r+1}) \leq (3/4) \sum_{i=1}^{m+r} \mu(s_i)$  can be re-written as:

$$\begin{aligned} \mu(s_{m+r+1}) &\leq \frac{3}{4} \sum_{i=1}^m \mu(s_i) + \frac{3}{4} \sum_{i=m+1}^{m+r} \mu(s_i) + \frac{3}{4} r \sum_{i=1}^m \mu(s_i) - \frac{3}{4} r \sum_{i=1}^m \mu(s_i) \\ &\quad + \frac{1}{4} \sum_{i=m+1}^{m+r} (3(m+r-i)-1) \mu(s_i) - \frac{1}{4} \sum_{i=m+1}^{m+r} (3(m+r-i)-1) \mu(s_i) \\ &= \frac{3}{4}(r+1) \sum_{i=1}^m \mu(s_i) + \frac{1}{4} \sum_{i=m+1}^{m+r} (3(m+r-i)-1) \mu(s_i) \\ &\quad - \frac{3}{4} r \sum_{i=1}^m \mu(s_i) + \frac{1}{4} \sum_{i=m+1}^{m+r} (3(1-m-r+i)+1) \mu(s_i) \\ &= \frac{3}{4}(r+1) \sum_{i=1}^m \mu(s_i) + \frac{1}{4} \sum_{i=m+1}^{m+r} (3(m+r-i)-1) \mu(s_i) + \delta \end{aligned}$$

where  $\delta = \frac{1}{4} \sum_{i=m+1}^{m+r} (3(1-m-r+i)+1)\mu(s_i) - \frac{3}{4}r \sum_{i=1}^m \mu(s_i)$ . If  $r = 1$ , then  $\delta = \mu(s_{m+1}) - \frac{3}{4} \sum_{i=1}^m \mu(s_i) \leq 0$  holds and the condition for the first-mover to play strategy  $s_m$  in Block 1 is satisfied. Now suppose that  $r \geq 2$ . We can rewrite  $\delta$  as follows:

$$\begin{aligned} \delta &= \frac{1}{4} \sum_{i=m+1}^{m+r} (3(1-m-r+i)+1)\mu(s_i) - \frac{3}{4}r \sum_{i=1}^m \mu(s_i) \\ &= \mu(s_{m+r}) + \frac{1}{4}\mu(s_{m+r-1}) + \frac{1}{4} \sum_{i=m+1}^{m+r-2} (3(1-m-r+i)+1)\mu(s_i) - \frac{3}{4}r \sum_{i=1}^m \mu(s_i) \\ &= \mu(s_{m+r}) + \frac{1}{4}\mu(s_{m+r-1}) + \gamma - \frac{3}{4}r \sum_{i=1}^m \mu(s_i), \end{aligned}$$

where  $\gamma = \frac{1}{4} \sum_{i=m+1}^{m+r-2} (3(1-m-r+i)+1)\mu(s_i)$ . Note that if  $r \geq 2$ , then  $\gamma < 0$ . Therefore, we have the following:

$$\begin{aligned} \delta &< \mu(s_{m+r}) + \frac{1}{4}\mu(s_{m+r-1}) - \frac{3}{4}r \sum_{i=1}^m \mu(s_i) \\ &\leq \left(\frac{3}{4}\right)^r + \frac{1}{4}\left(\frac{3}{4}\right)^{r-1} - \frac{3}{4}r \sum_{i=1}^m \mu(s_i) \\ &= \frac{3}{4}\left(\left(\frac{3}{4}\right)^{r-2} - r\right) \sum_{i=1}^m \mu(s_i). \end{aligned}$$

$r \geq 2$  implies that  $\left(\frac{3}{4}\right)^{r-2} - r < 0$ , so  $\delta < 0$  holds and the condition for the first-mover to play strategy  $s_m$  in Block 1 is satisfied.

Given that the second-mover always defected before her opponent by round  $n$  of Block 1 supergames, where  $m < n$ , we have  $\tilde{\mu}(s_k) = 0$  for all  $k \leq m$ . Therefore,  $\tilde{p}_l = 0$  holds for all  $l \leq m$ . By Proposition 1 it follows that the first-mover's Block 2 strategy is  $s_{m+t}$  for some  $t \geq 1$ .

#### Proof of Proposition 3

(a) By Proposition 1, the first-mover's Block 1 strategy is  $s_n$  if and only if  $\mu$  is such that  $p_k \geq \frac{4}{\sum_{i=n+1}^k (3\prod_{j=i}^{k-1} p_j) + 7\prod_{i=n}^{k-1} p_i}$  holds for all  $k \in \{n, \dots, 10\}$ . By similar logic to the proof of Proposition 2, if the first-mover's prior beliefs satisfy  $\mu(s_{k+1}) \leq (3/4) \sum_{i=1}^k \mu(s_i)$  for all  $k \in \{n, \dots, 10\}$  then the condition for the first-mover to play strategy  $s_n$  in Block 1 is satisfied. Given that the second-mover never defected before her opponent in rounds  $10, \dots, m$  of Block 1 supergames, where  $m < n$ , we have  $\tilde{\mu}(s_k) = 0$  for all  $k > m$ . Therefore,  $\tilde{p}_l = 1$  holds for all  $l \geq m$ . By Proposition 1 it follows that the first-mover's Block 2 strategy is  $s_{n-t}$  for some  $t \geq 1$ .

(b)  $\mu(s_{11}) > 3/7$  implies that  $\mu(s_{11}) > (3/4) \sum_{i=1}^{10} \mu(s_i)$  holds, so the condition for the first-mover to play strategy  $s_{11}$  in Block 1 is satisfied. Given that the second-mover never defected before her opponent in rounds  $10, \dots, m$  of Block 1 supergames, where  $m \leq 10$ , we have  $\tilde{\mu}(s_k) = 0$  for all  $k > m$ . Therefore,  $\tilde{p}_l = 1$  holds for all  $l \geq m$ . By Proposition 1 it follows that the first-mover's Block 2 strategy is  $s_{11-t}$  for some  $t \geq 1$ .

#### REFERENCES

- Ambrus A, Pathak PA (2011) Cooperation over finite horizons: A theory and experiments. *Journal of Public Economics* 95:500–512
- Andreoni J (1988) Why free ride? strategies and learning in public goods experiments. *Journal of Public Economics* 37:291–304
- Andreoni J, Croson R (1998) Partners vs. strangers: Random rematching in public goods experiments, working Paper
- Andreoni J, Miller JH (1993) Rational cooperation in the finitely repeated prisoners' dilemma: Experimental evidence. *Economic Journal* 103(418):570–585
- Blanco M, Engelmann D, Koch AK, Normann HT (2011) Preferences and beliefs in a sequential social dilemma: A within-subjects analysis, working Paper
- Bolle F, Ockenfels P (1990) Prisoners' dilemma as a game with incomplete information. *Journal of Economic Psychology* 11:69–84
- Bolton GE, Katok E, Ockenfels A (2005) Cooperation among strangers with limited information about reputation. *Journal of Public Economics* 89:1457–1468

- Brandts J, Charness G (2000) Hot vs. cold: Sequential responses and preference stability in experimental games. *Experimental Economics* 2:227–238
- Brandts J, Charness G (2011) The strategy versus the direct-response method: a first survey of experimental comparisons. *Experimental Economics* 14:375–398
- Brandts J, Figueras N (2003) An exploration of reputation formation in experimental games. *Journal of Economic Behavior and Organization* 50(1):89–115
- Bruttel LV, Güth W, Kamecke U (2012) Finitely repeated prisoners' dilemma experiments without a commonly known end. *International Journal of Game Theory* 41(1):23–47
- Camera G, Casari M (2009) Cooperation among strangers under the shadow of the future. *American Economic Review* 99:979–1005
- Camerer C, Weigelt K (1988) Experimental tests of a sequential equilibrium reputation model. *Econometrica* 56(1):1–36
- Clark K, Sefton M (2001) The sequential prisoner's dilemma: Evidence on reciprocity. *The Economic Journal* 111:51–68
- Cooper R, DeJong D, Forsythe R, Ross TW (1996) Cooperation without reputation: Experimental evidence from prisoner's dilemma games. *Games and Economic Behavior* 12:187–218
- Crawford VP, Costa-Gomes MA, Iriberrri N (2013) Structural models of nonequilibrium strategic thinking: Theory, evidence, and applications. *Journal of Economic Literature* 51(1):5–62
- Dal Bó P (2005) Cooperation under the shadow of the future: Experimental evidence from infinitely repeated games. *American Economic Review* 95:1591–1604
- Dal Bó P, Frechette G (2011) The evolution of cooperation in repeated games: Experimental evidence. *American Economic Review* 101:411–429
- Duffy J, Ochs J (2009) Cooperative behavior and the frequency of social interaction. *Games and Economic Behavior* 66:758–812
- Durkalski V, Palesch Y, Lipsitz S, Rust P (2003) Analysis of clustered matched-pair data. *Statistics in Medicine* 22:2417–2428
- Embrey M, Fréchette GR, Yuksel S (2014) Cooperation in the finitely repeated prisoner's dilemma, Working Paper
- Gächter S, Thoni C (2005) Social learning and voluntary cooperation among like-minded people. *Journal of the European Economic Association* 3(2-3):303–314
- Gong B, Yang CL (2010) Reputation and cooperation: an experiment on prisoner's dilemma with second-order information, working Paper
- Healy PJ (2007) Group reputations, stereotypes, and cooperation in a repeated labor market. *American Economic Review* 97(5):1751–1773
- Ho TH, Su X (2009) Peer-induced fairness in games. *American Economic Review* 99:2022–2049
- Irlenbusch B, Sliwka D (2005) Incentives, decision frames, and motivation crowding out – an experimental investigation, IZA Discussion Paper No. 1758
- Jung YJ, Kagel JH, Levin D (1994) On the existence of predatory pricing: An experimental study of reputation and entry deterrence in the chain-store game. *RAND Journal of Economics* 25(1):72–93
- Kagel JH, McGee P (2014) Team versus individual play in finitely repeated prisoner dilemma games, Working Paper
- Kahn LM, Murnighan JK (1993) Conjecture, uncertainty, and cooperation in prisoner's dilemma games. *Journal of Economic Behavior and Organization* 22:91–117
- King-Casas B, Tomlin D, Anen C, Camerer CF, Quartz SR, Montague PR (2005) Getting to know you: Reputation and trust in a two-person economic exchange. *Science* 308:78–83
- Kreps DM, Milgrom P, Roberts J, Wilson R (1982) Rational cooperation in the finitely repeated prisoners' dilemma. *Journal of Economic Theory* 27:245–252
- Larocque D (2005) Statistical Modeling and Analysis for Complex Data Problems, chap The Wilcoxon Signed-Rank Test for Cluster Correlated Data, pp 309–323
- Murnighan JK, Roth AE (1983) Expecting continued play in prisoner's dilemma games. *Journal of Conflict Resolution* 27:279–300
- Neral J, Ochs J (1992) The sequential equilibrium theory of reputation building: A further test. *Econometrica* 60(5):1151–1169
- Normann HT, Wallace B (2012) The impact of the termination rule on cooperation in a prisoner's dilemma experiment. *International Journal of Game Theory* 41(3):707–718
- Nowak MA, Sigmund K (2005) Evolution of indirect reciprocity. *Nature* 437:1291–1298
- Radner R (1986) Can bounded rationality resolve the prisoners' dilemma? In: Mas-Colell A, Hildenbrand W (eds) *Contributions to Mathematical Economics*, North-Holland, Amsterdam, pp 387–399
- Reuben E, Suetens S (2012) Revisiting strategic versus non-strategic cooperation. *Experimental Economics* 15:24–43

- Roe BE, Wu SY (2009) Do the selfish mimic cooperators? experimental evidence from finitely-repeated labor markets, working Paper
- Roth AE, Murnighan JK (1978) Equilibrium behavior and repeated play of the prisoner's dilemma. *Journal of Mathematical Psychology* 11:189–198
- Samuelson L (1987) A note on uncertainty and cooperation in a finitely repeated prisoners' dilemma. *International Journal of Game Theory* 16:187–195
- Schwartz ST, Young RA, Zvinakis K (2000) Reputation without repeated interaction: A role for public disclosures. *Review of Accounting Studies* 5:351–375
- Selten R, Stoecker R (1986) End behavior in sequences of finite repeated prisoner's dilemma supergames: a learning theory approach. *Journal of Economic Behavior and Organization* 7:47–70
- Tingley DH, Walter BF (2011) The effect of repeated play on reputation building: An experimental approach. *International Organization* 65:343–365