# Using Team Discussions to Understand Behavior in Indefinitely Repeated Prisoner's Dilemma Games[†]

*By* DAVID J. COOPER AND JOHN H. KAGEL*

*We compare behavior of two person teams with individuals in indefinitely repeated prisoner dilemma games with perfect monitoring. Team discussions are used to understand the rationale underlying these choices and how these choices come about. There are three main findings: (i) Teams learned to cooperate faster than individuals, and cooperation was more stable for teams. (ii) Strategies identified from team dialogues differ from those identified by the Strategy Frequency Estimation Method. This reflects the improvisational nature of teams' decision making. (iii) Increasing cooperation was primarily driven by teams unilaterally cooperating in the hope of inducing their opponent to cooperate. (JEL C72, C73, C92)*

This paper reports an indefinitely repeated prisoner's dilemma (IRPD) game experiment where the decision making "agents" were either freely interacting two-person teams or individuals. The primary goals of the paper are twofold. First, to compare the behavior of teams and individuals in IRPD games with perfect monitoring where mutual cooperation is both consistent with equilibrium and risk dominant. Many strategic interactions in economics involve groups of individuals. If the behavior of groups differs substantially from individuals, conclusions based solely on the decisions of individuals will not be broadly applicable. Second, to use team dialogues to understand what motivates agents' choices and how these come about. Economists have become increasingly interested in using process data

(e.g., response times, eye tracking, and fMRI) to understand decision making.[1] Analyzing the content of team dialogues provides direct insights into the thought processes underlying agents' choices that cannot be easily obtained using other methods. This makes it possible to directly observe *how* and *why* strategies came about. Given that the broad patterns of play were similar for teams and individuals, team discussions presumably offer insights into the motivation underlying individuals' choices as well.

Comparing teams and individuals, cooperation was rare for both in early supergames, with teams cooperating somewhat less than individuals. Cooperation increased significantly faster *across* supergames for teams than individuals, so that in later supergames teams cooperated at significantly higher rates. *Within* supergames, play was significantly more stable for teams than individuals. Play by teams was also more stable *between* supergames; once a team switched to a cooperative strategy, they rarely switched back.

Teams rarely discussed strategies as a game theorist would. Rather than pre-specifying a full, state contingent plan, strategies were typically incompletely specified, with teams improvising in response to unanticipated choices by their opponents. Despite this, we can almost always identify a strategy that corresponds to their initial plan for making choices within a supergame.

The modal strategy identified from team dialogues in early supergames was Always Defect (AD). Over half of all teams who started out with AD switched to more cooperative strategies, primarily Grim Trigger and its lenient variants, so that in later supergames a majority of teams employed cooperative strategies. Returning to AD, even temporarily, was rare.

Comparing the distribution of strategies identified from team dialogues with estimates based on the Strategy Frequency Estimation Method (SFEM; Dal Bó and Fréchette 2011), SFEM puts more weight on variants of Tit-for-Tat (TFT, STFT, etc.) than variants of Grim Trigger, while the coding based on team dialogues does the opposite. This reflects the improvisational nature of teams' strategies. Teams did not try to anticipate all possible histories of play. Instead, they used simple plans covering the first few stage games (e.g., "My plan is choose C [cooperate] first for a few times, and if the other team keeps choosing D [defect], we will switch to D.") and then adjusted as needed. Choice patterns that appeared like Tit-for-Tat often reflected improvised reactions to their opponent, attempting to coordinate on mutual cooperation within a supergame, rather than a preconceived strategy.

Analysis of team dialogues makes it possible to identify the rationales underlying teams' adoption of strategies. Their frequent initial use of AD was primarily motivated by fear of their opponents defecting and a desire for the safety that AD provides. Teams typically had long conversations prior to first trying a cooperative strategy, with discussions stretching across multiple stage games and even across supergames. Realizing that *mutual* cooperation pays more in the long run, they often tried cooperating for a few stage games to see if the other team would go along. They discussed this in terms of "leading by example," hoping that their actions

---

[1] See Fehr and Glimcher (2013) and Cooper, Krajbich, and Noussair (2019) for recent surveys.

would send a message when direct communication was not possible. As one team noted, "this is all so hard without communication" "I know if we could just send them like one sentence we'd have it made." Teams varied in the specifics of how they conceptualized and executed this approach, making possible the identification of different strategies, but ultimately these were slight variations on the same basic theme.

Team dialogues also explain the stability of team play relative to individuals. When teammates disagreed about whether to stick with the status quo or make a change, the status quo usually won (87 percent). Teammates provided a check on switching that did not exist for individuals.

The comparison of teams and individuals raises a question whether behavior differs due to the presence of a teammate per se, or because of joint decision making and the associated communication between teammates. To distinguish between these two possibilities, a silent partners treatment was implemented. Subjects played the game in fixed pairs, like the team treatment, but one teammate was solely responsible for making decisions with no input from their silent partner. Behavior in the silent partner treatment differs little from the individual treatment, indicating that the effect of team play was largely due to a combination of bilateral communication and joint decision making.

The design and procedures used here parallel those employed in Kagel and McGee (2016) to compare individual and team play in finitely repeated prisoner's dilemma (FRPD) games. This makes possible a clean comparison of the differences between teams and individuals in IRPD and FRPD games. There are strong similarities between the two with inexperienced teams less cooperative than individuals in early stages of early supergames, only to become more cooperative with experience in early stage games. Likewise, play by teams was more stable than individuals. The differences between teams and individuals likely reflect underlying differences, present in many settings, between how teams and individuals make choices.

It is inherently interesting to understand *how* and *why* teams chose strategies, but this understanding also has implications beyond the experiment reported here. Teams approached IRPD games with a combination of simplicity and sophistication. Rather than fixing a detailed plan in advance, they typically start with a simple, incomplete, initial plan followed by improvisation in response to the behavior of their opponent. At the same time, leading by example showed a sophisticated ability to anticipate and manipulate other agents' behavior. The simple, flexible approach taken by teams is portable. Rather than only applying to the environment studied here, it is an approach that seems easily applied to other repeated games. The sophistication exhibited by teams and their tendency to improvise has important implications for how learning in IRPD games should be modeled, an issue discussed at length in the conclusion.

In terms of technique, the approach employed here for coding team dialogues is a departure from how economists have previously coded communication. Most previous coding exercises have been done at a granular level, but this would have missed vital content from conversations that extended over multiple stage games or even supergames. Identifying what strategies were being used (and why) required synthesizing the content of long-running conversations, in conjunction with choices, rather

than focusing on what was said at any single point in time. These procedures are likely to be applicable to other experimental settings as well.

The outline of the paper is as follows. Section I briefly reviews past research with IRPD games. Section II reports the experimental design and procedures, and Section III lays out hypotheses for the data based on the relevant theory and the preceding literature. Section IV compares differences between teams and individuals for IRPD games. Section V describes the main features of the team dialogues, explains how strategies are identified from team discussions, investigates how strategies evolve with experience, compares the strategies identified from team dialogues with the distribution of strategies estimated by SFEM, and examines the rationales underlying teams' choices. Section VI contrasts the results reported here with results from previously reported FRPD games. Section VII reports results from the silent partner treatment. Section VIII summarizes the results and discusses the implications of what has been reported.

## I. Prior Research

There have been numerous experimental studies of IRPD games. Dal Bó and Fréchette (2018) provide an extensive survey of the experimental literature. This paper departs from the existing literature along two important dimensions. First, most of the economics research on IRPD games involves individuals making decisions. The focus here is on teams, specifically on the processes and motivation for their choices. Cason and Mui (2019) is a notable exception to the focus on individuals, comparing play in three-person teams with individuals. Most of their analysis considers games with imperfect monitoring, with data primarily based on subjects' choices over a menu with predefined strategies (e.g., tit-for-tat or grim trigger). In a control treatment with perfect monitoring (like this experiment), *both* teams and individuals achieved very high cooperation rates (80–90 percent). These high cooperation rates reflect payoff matrices and continuation probabilities designed to achieve reasonable cooperation rates with imperfect monitoring, but leave little room for distinguishing between individuals and teams with perfect monitoring.[2] Teams in Cason and Mui could discuss their strategy choices. The paper describes few details of their formal coding of these discussions, which only play a minor role in the paper, but they do make use of samples of team discussions to illustrate features of team play. Identifying strategies from chat is a non-issue due to their use of direct elicitation.[3] The focus of Cason and Mui, compared to the present paper, is quite different. More than anything else, the present paper is about what can be learned from the team chats. Cason and Mui is largely concerned with what strategies teams chose compared to individuals.

The second important difference between this paper and the existing literature is the use of team dialogues to identify what strategies are used and why. The most

---

[2] They have a substantially smaller basin of attraction for Always Defect (SizeBAD): 0.05 compared to 0.26 in our design. See Section III for a definition and discussion of SizeBAD in relation to cooperation rates.

[3] They report relationships between appeals to game theoretic reasoning and choice of AD and between concerns about the effect of noisy implementation of chosen actions and the use of lenient strategies.

common method for identifying strategies is the Strategy Frequency Estimation Method (SFEM) introduced in Dal Bó and Fréchette (2011). This technique uses maximum likelihood estimation of a mixture model to estimate the distribution of strategies across the population. There are also a small number of papers that use direct elicitation, asking agents to specify a full strategy, either directly or by making a choice from a menu, at the beginning of each supergame (Romero and Rosokha 2018, 2019; Dal Bó and Fréchette 2019; Cason and Mui 2019).

Each method of identifying strategies in IRPD games has its strengths and weaknesses, and which is best depends on the purpose of the exercise. Econometric approaches like SFEM are unintrusive, but only indirectly measure what strategies were used based on a predefined set of strategies. Such methods are best for estimating strategies employed when the researchers run a standard experiment with individuals making direct choices. This makes sense when identification of individual agents' strategies is not the primary purpose of the paper. Direct elicitation, by its very nature perfectly identifies what strategies were used, but is the most intrusive method. Direct elicitation often limits subjects to a prespecified list of strategies. Even in cases where any strategy can be constructed, the exercise implicitly directs subjects toward completely specified strategies rather than the incomplete strategies identified in our experiment.[4] Direct elicitation is the best approach if the goal is to unambiguously identify strategies. Analysis of team chat lies somewhere between econometric methods and direct elicitation. Making decisions in a team obviously affects how choices are made, but discussions between teammates are a natural part of this process and recording the chats is nonintrusive. Identification of strategies from team chat is more direct than econometric methods, but does involve judgment by the coders, making it less accurate than direct elicitation. The true advantage of analyzing team dialogues is the ability to understand how and why decisions were made. If the goal is to better understand the nature of strategies and the rationale underlying agents' strategies, chat analysis is likely to be the best approach.

There is an important line of research in the social psychology literature concerned with differences between individuals and teams in PD games. These experiments are quite different in structure from those commonly employed in economics. They typically involve a single supergame where agents are told they will be paired with the same opponent for between $t$ and $t + n$ stage games, with the actual stopping point occurring at an unknown point in that interval. The key finding from this literature is that teams are less cooperative than individuals (referred to as the "discontinuity effect"). This is attributed to greater fear and greed on the part of groups than individuals (see Wildschut et al. 2003 and Wildschut and Insko 2007 for surveys).

Beyond IRPD games, there is a growing literature in experimental economics comparing how individuals and teams behave in games (e.g., Cooper and Kagel 2005; Kocher and Sutter 2005; Feri, Irlenbusch, and Sutter 2010; Maciejovsky et al.

---

[4]Romero and Rosokha (2018, 2019) allowed subjects to construct strategies in an almost unlimited fashion. In the latter paper, subjects could also revise their strategies midgame. In both cases, strategies are complete plans for how the game is to be played. Even if this plan can be revised, having subjects specify complete plans presumably has an effect on how subjects conceive of strategies. Dal Bó and Fréchette (2019) included a phase in which subjects need not follow their chosen strategies. The same concern about framing applies.

Table 1. Stage Game Payoffs

|   | A | B |
|---|---|---|
| A | 105<br>105 | 5<br>175 |
| B | 175<br>5 | 75<br>75 |

2013; Casari, Zhang and Jackson, 2016). The typical finding is that teams are more "rational" than individuals; specifically, teams are more likely to use a theoretically optimal strategy and are faster to learn how to maximize their payoffs. There is no optimal strategy in IRPD games, but, given that cooperative strategies earned higher payoffs than AD, finding that teams switched to cooperative strategies faster than individuals is consistent with this literature.

## II. Experimental Design and Procedures

Throughout this paper "agent" is used as a generic term for either individuals or two-person teams. Agents played a simultaneous move, indefinitely repeated prisoner's dilemma (IRPD) game with perfect monitoring using the stage game payoffs reported in Table 1 (own payoffs are in red, opponent's payoffs are in blue).

The continuation probability was $\delta = 0.90$, yielding an expected supergame length of ten stage games. Following each supergame, agents were randomly rematched with the restriction that no two agents were matched in consecutive supergames. The instructions stressed that at the end of each stage game "there is a 90 percent chance of another round [stage game] for that match [supergame] and a 10 percent chance you will move on to another match with another team/individual." After the last stage game within a supergame, agents were notified that their match had ended and that they would be starting another match with a different (randomly chosen) agent. Judging from the team chats, agents had no difficulty telling when they were starting a new supergame with a different opponent.

Using a between-subjects design, all agents in a session were either individuals or teams with no mixing between the two. There were six individual agent sessions with 14–18 subjects in each session (104 total subjects). Sessions lasted 90 minutes, conducting as many supergames as possible within the allotted time. Four sessions had 13 supergames, the other two had 12. Individuals had up to one minute to make their choice, but this was never a binding constraint.

In the team treatment, subjects were randomly matched at the beginning of the session. Teammates remained the same throughout the session. Instructions were essentially the same as for individual sessions (see online Appendix A), except that teams were told to "make decisions jointly." Each teammate could enter a choice, with the team's choice implemented if choices agreed and did not change for five seconds. In the first two stage games of each supergame, teammates had two minutes to discuss their choices and reach an agreement. This was reduced to 40 seconds for all subsequent stage games. The default options, in case teammates could not agree, are

reported in the instructions (see online Appendix A). Teammates reached an agreement on what action to take in the overwhelming majority of cases ($>$ 99 percent). To facilitate reaching agreement, teammates could send messages back and forth. They were told to use the "chat box to discuss your choices, and come to an agreement regarding what choice to make." Subjects were instructed to be civil to each other, not use profanity, and to not identify themselves. The instructions stressed that other teams could not see their messages.

There were 6 team sessions with between 8 and 12 teams in each session for a total of 58 teams (116 subjects). The first two sessions had only six and seven supergames respectively; these were scheduled to last two hours as the amount of time the team chats would take was not anticipated. Subsequently, session time was increased to two and a half hours, along with modest reductions in the time teammates could discuss their choices.[5] The remaining four team sessions had between nine and twelve supergames.

The six individual and six team sessions were paired, using matching seeds to generate the random length of the supergames (see online Appendix B for a session list). This guaranteed parallelism between the two treatments in terms of the length of supergames. The main reason for variability in the number of supergames across sessions, beyond the timing issue mentioned above, is that the different seeds led to different lengths of supergames.

Payoffs were denominated in experimental currency units (ECUs) which were converted into dollars at the rate of $1 = 250$ ECUs. Payoffs were summed over all stage games of all supergames, converted to dollars, and paid in cash at the end of an experimental session. Each member of a team received the team's full payoff. Earnings averaged $44.98 per subject, with a typical session lasting 90 minutes for individuals and 135–150 minutes for teams. Subjects were recruited from the Ohio State University undergraduate population using ORSEE (Greiner 2015). The software was programmed using zTree (Fischbacher 2007).

## III. Predictions

With sufficiently patient agents, IRPD games with perfect monitoring have an infinite set of subgame perfect equilibrium outcomes.[6] Cooperative outcomes are consistent with equilibrium, but noncooperative play is an equilibrium as well. A number of criteria have been explored to better predict when cooperation is more likely, the most popular of which reduces the game to a normal form game with only grim trigger (Grim) and always defect (AD) available as strategies. *SizeBAD* (Dal Bó and Fréchette 2011) is defined as the size of the basin of attraction for AD. *SizeBad* measures the dynamic stability of cooperation. If *SizeBad* $< 1$, mutual play of Grim is a subgame perfect equilibrium, and if *SizeBad* $< 0.50$, mutual play of Grim is risk dominant (Harsanyi and Selten 1988). Empirically, cooperation rates

---

[5] In the first two sessions, teams had two and a half minutes to discuss their choices in the first two-stage games, with one minute after that. Subject feedback indicated that shorter times were more than adequate to reach agreement, so teams in the remaining sessions were given two minutes to discuss their choices in the first two-stage games, with 40 seconds after that.

[6] Friedman (1971); Aumann and Shapley (1994); and Fudenberg and Maskin (1986).

are a decreasing function of *SizeBad* (Dal Bó and Fréchette 2011 and 2018).[7] The combination of stage game payoffs and continuation probability used here yields *SizeBAD* = 0.26; cooperation is risk dominant, but not so strong that we'd expect universal cooperation.

As noted previously, the experimental literature finds that teams are more likely than individuals to use a theoretically optimal strategy in a game, and are faster to learn how to maximize their payoffs. The former is of no help in making predictions, since the benefits of specific strategies depend on what equilibrium is played, but the latter implies that mutual cooperation rates will increase faster for teams. The discontinuity effect documented in the psychology literature, as described above, suggests that cooperation rates will initially be lower for teams.

Finally, previous experiments find a strong relationship between response times and cooperation in one-shot games (Rand, Greene, and Nowak 2012), with more cooperative individuals having slower response times. This is attributed to cooperative play requiring greater deliberation, and remains true when individuals are prompted to respond either more or less rapidly. Team play forces slower decisions due to the need to discuss the team's choices. As such, teams might be expected to be more cooperative. This finding refers to the level of cooperation and does not deal with whether cooperation will grow more rapidly for teams versus individuals.

## IV. Experimental Results, Individuals versus Teams

Holding the seed fixed, the number of supergames completed was always less in the team treatment. To maintain parallelism, the dataset is restricted to data prior to and including the *final common supergame*, defined as the last supergame played by *both* individuals and teams using the *same* seed. For example, if teams played 9 supergames with a given seed and individuals played 12 with the same seed, in both cases data used is from the first nine supergames.

The notation SGx refers to the x-th supergame in an experimental session (i.e., SG1 for the first supergame, SG2 for the second supergame, etc.) and Stx refers to the x-th stage game within a supergame (i.e., St1 for the first stage game, St2 for the second stage game, etc.). Unless stated otherwise, statements about statistical significance comparing teams and individuals are based on Wilcoxon matched-pairs signed-rank tests, where observations are session averages paired by seed class. Likewise, statements about the statistical significance of changes over time within a session (e.g., does mutual cooperation differ between the first and last supergame) are also based on Wilcoxon matched-pairs signed-rank tests. These are weak tests that are biased in favor of Type II errors. Regressions yielding similar results are reported in online Appendix C.

Throughout this section, an observation is defined as a single play of a stage game. There are three possible outcomes for each observation: mutual cooperation (CC), mutual defection (DD), and mixed (CD). Mutual cooperation is the primary measure of cooperation used throughout the paper. Stressing mutual outcomes

---

[7] Also see Lugovskyy, Puzzello, and Walker (2018) and Blonski, Ockenfels, and Spagnolo (2011).

rather than individual agents' choices is largely a matter of convenience; individual cooperation and mutual cooperation are highly correlated, and it is redundant to describe results for both. Analysis reported in online Appendix C shows that the main conclusions are not affected by using mutual cooperation as the measure of cooperation rather than cooperation by individual agents.[8]

Figure 1A reports the rate of mutual cooperation over the first ten stage games (St1−St10). The data are broken down by individuals or teams, and early (SGs 1−3) or late (SG $\geq$ 4) supergames. As another way of seeing how the data changes with experience, Figure 1B shows the average mutual cooperation rate in St1 of SGs 1 − 8; after this point in time, more than half of all agents had dropped out of the sample due to sessions ending. For parallel figures based on individual cooperation rates, along with a brief discussion, see online Appendix C.

Looking at Figure 1A, notice that the frequency of mutual cooperation changed little after the first few stage games.[9] Differences between treatments were largely driven by what happened in St1. Given this, the following analysis focuses on mutual cooperation *in* St1. This measure has the advantage of not being affected by the differing length of supergames and is highly correlated with mutual cooperation in later stage games ($\rho = 0.69$). See online Appendix C for results showing that our conclusions are not affected by using data from all stage games.

Mutual cooperation in St1 was *lower* for teams than individuals in SG1 (10.3 percent versus 19.2 percent), but with experience teams overtook individuals: by the final common supergame, mutual cooperation in St1 was *higher* for teams (55.2 percent versus 36.5 percent). Comparing SG1 and the last common supergame, mutual cooperation in St1 increased significantly for teams ($n = 6$; $z = 2.20$; $p = 0.028$) but not for individuals ($n = 6$; $z = 1.05$; $p = 0.292$), with the increase significantly larger for teams compared to individuals ($n = 12$; $z = 1.78$; $p = 0.075$). Agents who cooperated in St1 earned more than those who did not.[10] The faster increase in mutual cooperation on the part of teams is consistent with past findings in other settings (cited in Section I) that teams learn to maximize payoffs more rapidly than individuals.

OBSERVATION 1: *Mutual cooperation increased faster with experience for teams than for individuals.*

While *average* levels of cooperation were stable across stage games, this hides a fair amount of switching between outcomes (mutual cooperation, mutual defection, or mixed) within *individual* supergames. A "switch" occurs when the outcome for the current stage game differs from the outcome in the previous stage game within

---

[8] In 89 percent of all stage games, agents either mutually cooperated or mutually defected.

[9] The early increase in mutual cooperation was not statistically significant based on a probit regression using data from all stage games. The dependent variable is a dummy for mutual cooperation, and independent variables include controls for the treatment, supergame, and seed class. A dummy for stage games greater than or equal to 3 has a positive parameter estimate, but is not significant (est. $= 0.034$; SE $= 0.028$; $p = 0.223$).

[10] The difference in average payoffs per stage game between those agents who cooperate in St1 versus those who don't was 90 versus 85 ECUs for individuals and 91 versus 83 for teams.
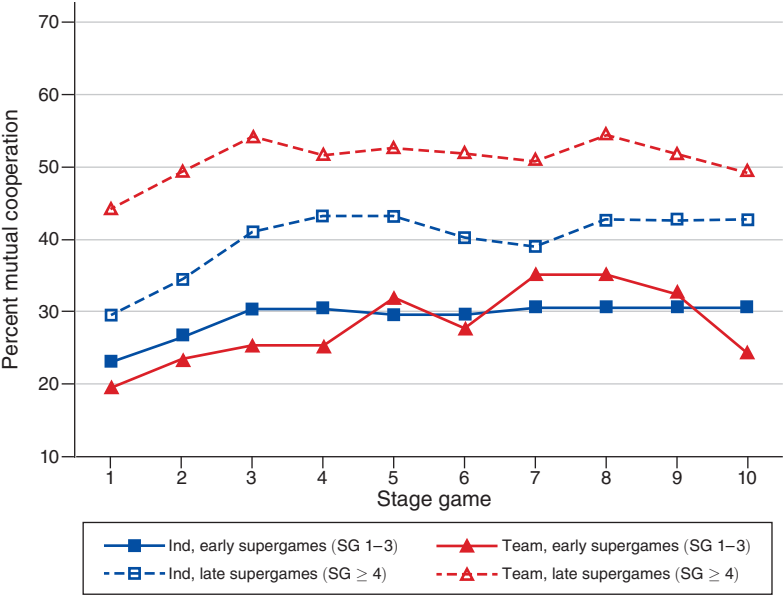
FIGURE 1A. MUTUAL COOPERATION ACROSS STAGE GAMES

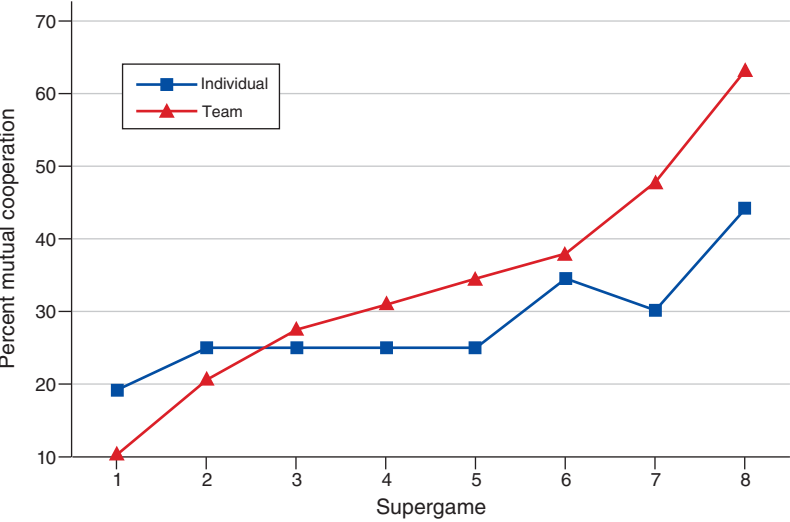*Note:* Figure 1A is based on 9,498 observations.



FIGURE 1B. MUTUAL COOPERATION IN STAGE GAME 1 ACROSS SUPERGAMES

*Note:* Figure 1B is based on 723 observations.

TABLE 2. NUMBER OF SWITCHES PER SUPERGAME

|  |  | Individual | Team |
|---|---|---|---|
| All observations | Average | 1.43 | 0.85 |
|  | Number of observations | 344 | 194 |
| Mutual cooperation (CC) | Average | 0.93 | 0.69 |
|  | Number of observations | 100 | 70 |
| Mutual defection (DD) | Average | 1.48 | 0.44 |
|  | Number of observations | 79 | 52 |
| Mixed (CD) | Average | 1.71 | 1.31 |
|  | Number of observations | 165 | 72 |

a given supergame.[11] Table 2 reports the average number of switches per supergame, excluding very short supergames with only one or two stage games (top row). Results are also reported separately for each possible outcome in St1 (rows 2–4, respectively).

The number of switches was significantly lower for teams than individuals across all observations ($n = 6$; $z = 1.99$; $p = 0.046$). Teams had fewer switches than individuals for all initial outcomes, but the difference was largest following mutual defection in St1. The difference was slightly smaller in late supergames (SG $\geq 4$), but the average number of switches remained significantly higher for individuals (1.26 versus 0.76; $n = 6$; $z = 1.99$; $p = 0.046$).

The likelihood of a switch was higher for individuals following either mutual cooperation (3.9 percent versus 2.5 percent) or mutual defection (5.0 percent versus 1.7 percent). These small differences in the probability of switching had substantial cumulative effects relative to outcomes in St1. This is illustrated in Figure 2, which shows the fraction of observations where the outcome *differed* from the *initial* outcome in St1, distinguishing between starting with mutual cooperation (left panel) and mutual defection (right panel). For example, consider *pairs* of opponents that mutually cooperated in St1. In the team treatment, 0 percent switched to a *different* outcome in St2, compared to 8.9 percent for individuals. For St3, only 1.4 percent of teams are no longer mutually cooperating, far lower than the 16.0 percent figure for individuals. Greater stability is a double-edged sword; teams are better at sustaining mutual cooperation than individuals, but worse at escaping from mutual defection.

The greater stability of team play applied *between* supergames as well. Classify an agent as having made a "switch between supergames" if their action in St1 of the current supergame differs from St1 in the previous supergame. The proportion of switches between supergames was 23 percent for individuals versus 14 percent for teams, which is a weakly significant difference ($n = 6$; $z = 1.78$; $p = 0.075$). Analysis of the team chat, reported below in Table 4, shows that teams tend to switch from Always Defect to cooperative strategies, but rarely switch back. This may be

---

[11] For example, suppose a pair of agents have the outcomes C/C, C/C, C/D, D/D, D/D, and C/C in a supergame lasting for six stage games. There are three switches, in St3, St4, and St6.
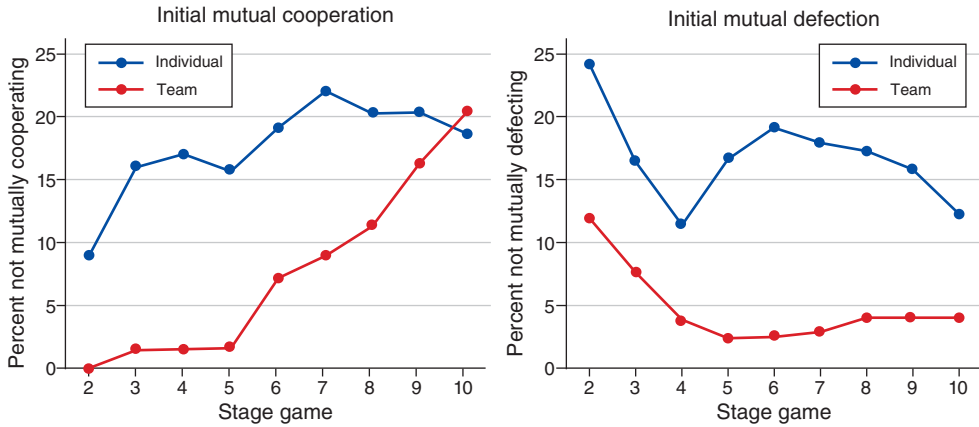
FIGURE 2. STABILITY, LIKELIHOOD OF NOT PLAYING INITIAL OUTCOME

*Note:* The left panel is based on 995 observations, and the right panel includes 1,272 observations.

less true for individuals than teams. Not only do teams switch between supergames less than individuals, they are also particularly unlikely to switch away from cooperation in St1. Teams are more than twice as likely to switch from defection to cooperation in St1 than vice versa (19 percent versus 9 percent), but this gap is smaller for individuals (25 percent versus 20 percent).

OBSERVATION 2: *Play was more stable for teams than individuals, both within supergames and between supergames.*

The nonparametric tests reported above are useful but conservative and somewhat limited, since there are no controls either for the varying length of supergames, or for agents' prior experience, which are known to affect cooperation rates. Online Appendix C reports regressions that control for these issues. The regressions confirm Observations 1 and 2, with statistical significance at the 5 percent level or better in both cases. Observations 1 and 2 continue to hold with a variety of changes in the specification (e.g., using all stage games rather than St1, using individual cooperation rather than mutual cooperation).

### V. Analysis of Team Discussions and Strategies

The main innovation of this paper is the use of teams' discussions to understand the thought processes and motivations underlying teams' "strategies": What strategies teams used, why did they use them, and how (and why) did they change with experience?

For all discussions, quotation marks separate the different team members' messages, with choices "A" and "B" changed to "C" and "D" to aid the reader. Spelling and grammatical errors are not corrected.

### A. *Features of Team Discussions*

Several common features of the team discussions dominate the analysis and interpretation of the data: (1) once a strategy was adopted, teams often continued to use it without restating it, although it is clear from their choices that they were following the same strategy; (2) ideas were often developed across multiple stage games, with choices frequently based on discussions in earlier stage games as well as earlier supergames; and (3) teams frequently did not specify a complete set of state contingent actions, improvising instead. The following examples illustrate these three features.

(1) Continuing to use a strategy without restating it: This team first used a cooperative strategy in SG8. They had started all prior supergames with D but had multiple discussions about possibly using a lenient Grim strategy. Their play from the beginning of SG8 was consistent with Grim with leniency for 3 stage games (Grim3), with a brief reiteration of the logic ("D?" "Let's try C a couple more times to see if they catch on"). For the subsequent four supergames (SG9–SG12), play was 100 percent consistent with Grim3 although they never discussed their strategy again. It was coded as such.

(2) Using a strategy based on earlier extended discussions: This team switched to Grim at the beginning of SG2. They had discussed this switch across multiple stage games in SG1 (e.g., "think C until opponent choose D is a good strategy to start each match?" "OK") but did not talk about their strategy in SG2. The relevant discussions had all taken place earlier.

(3) The improvisational nature of team decision making: The sequence of actions this team experienced was as follows, with own choice listed first: D/D, D/D, C/D, C/D, D/C, C/D, C/C, followed by continued mutual cooperation. They started out not trusting the other team, choosing D out of fear of getting the "sucker payoff." In St2, they discussed trying to achieve mutual cooperation and in St3 switched to C:

> St2: "I'm wondering if we can coordinate with them and pick C. Then we would both get 105 unless they choose D again, which they probably will... Lets do D again" "yes that's what I'm thinking..." "One of our teams needs to get smart and always choose C so we all get the most lol" "true! we need to coordinate."

> St3: "what about now..." "If one of us makes the jump to C, there is a chance that the other team will keep doing D. Is this a test of our selfishness? lol" "ok!" "Shall we try C?"

After trying to achieve mutual cooperation in St3 and St4, they gave up in St5.

St5: "Well that other team is selfish lol lets go with D. We tried" "i think no one will choose C … let's choose D" "Yeah. If we all take a chance and choose C, we will all get more!"

In St5, the other team switched to C just as they returned to D. In St6, they scrambled to recover from their poor timing, chosing C while the other team chose D.

St 6: "gosh dang it lol … they picked C … Back to C?" "let's do it lol."

St7: "well … C again?" "why other team just not coordinate with us... D maybe?" "just read our minds, jeez" "lol … I think we should stick with C one more time."

At this point, mutual cooperation was finally achieved. There was no over-arching plan behind this team's choices, so it does not fit neatly into the paradigm of picking a strategy at the beginning of a supergame, in the standard game theoretic sense. Rather they changed their mind about their plan in St2 and improvised after that.

It is worth noting that there is no evidence of *intentional* mixing in the team discussions as economists would define it. This contrasts with findings from Breitmoser (2015) and Romero and Rosokha (2021).

## B. *Coding Team Strategies*

To quantify the content of team discussions, two coding exercises were conducted. The first was a game-by-game coding of the dialogues in which the unit of coding was a stage game (specifically, a dialogue started after teammates learned the outcome of the previous stage game). The two authors went separately through the team chats from a randomly selected session and developed categories to be coded. These categories included discussing what action (C or D) to choose, reasons for choosing either an action or strategy, and generic questions about game play (e.g., asking a teammate about the continuation probability).

Coders were also asked to identify which of six simple strategies agents used: Always Defect (AD), Always Cooperate (AC), Grim Trigger (Grim), Tit-for-Tat (TFT), Suspicious Tit-for-Tat (STFT), and Win Stay, Lose Shift (WSLS).[12] Based on a meta-study of IRPD experiments, Dal Bó and Fréchette (2018) found that these generally accounted for more than three-quarters of the strategies used. Teams often discussed versions of Grim Trigger with leniency that permitted two or three defections before triggering punishment (Grim2 and Grim3), coded as Lenient Grim.[13]

---

[12] STFT is the same as TFT, except it starts with defect rather than cooperate. WSLS starts with cooperate, and then cooperates if and only if both players used the same action in the previous stage game.

[13] Grim2 and Grim3 were difficult to distinguish, hence,, we combined them into a single category. We also added more complex versions of TFT (e.g., tit for two tats; two tits for two tats etc.) to the coding scheme, paralleling the addition of more complex versions of Grim, but no examples were observed.

Complex versions of STFT (labeled "Complex STFT")[14] and Grim with Counting were added to the strategy set at the suggestion of the coders. Complex STFT refers to teams that started with D but quickly switched to C in an attempt to achieve cooperation (possibly in response to cooperation by the other team).[15] Grim with Counting describes cases where a team started with Grim but at some preset point, usually close to St10, unilaterally defected on the grounds that the supergame was likely to end (i.e., these teams suffered from the gambler's fallacy). Grim with Counting had received little attention in the existing literature but was previously reported in Romero and Rosokha (2018). The dialogue below is an example from a team that unilaterally defected in St9:

> St7: "when do you want to go D?" "I say 10 since they can last this long"
>
> St8: "10?" "9" "cool with me case thats the average amount"
>
> After defecting in St9…
>
> St10: "okay we gotta stick with D now lets hope it ends soon" "true"

Two graduate student RAs coded the dialogues independently. The coders were provided with copies of the coding categories along with explanations of each category.[16] The RAs were not told about any hypotheses the authors had, and conversations between the RAs and the coauthors were limited to clarifications of the coding scheme rather than suggestions about how any specific discussion should be coded. The coders were instructed that their task was to quantify the content of messages rather than interpreting the messages. The coders were free to code multiple categories for a single stage game. The coders were encouraged to add categories if they felt the researchers had missed something; complex STFT and Grim with Counting were both suggested by the coders. Online Appendix D reports the full set of coding categories, frequencies for each category, and Cohen's kappa for agreement rates between the two coders, which are generally quite high.

---

[14] STFT and Complex STFT are not identical. For instance, the following sequence: D/C, C/D, C/C is not consistent with STFT which would have had D/C in St3.

[15] Complex STFT often led to mutual cooperation. This might require multiple plays of C/D before abandoning C or settling down to mutual cooperation. For example, this team's dialogue prior to St2 stated the underlying strategy: "Should we switch to C… I feel like yes" "C?" "they'll go C if we do a couple times" "ok" "why not haha … worth a shot … maybe end up with 5 once but if we get 5 twice we will switch back to D." A second example of Complex STFT did not involve responding to initial cooperation by the other team. Play began D/D, and then one team unilaterally switched to C. Prior to St1, this team briefly stated their strategy: "D first then ride C" "I'd say so." They reiterate this strategy again before flipping to C in St2: "i think we might have to stick to D this round" "C for 2 times. If they don't switch, then we go back to D" "ok." The two teams achieved mutual cooperation in St3 and subsequently.

[16] For instance, Table D1 in online Appendix D lists "Myopia" as a category. The coders' instructions added the following explanation: Myopia has a number of possible characteristics all of which should lead to classification under myopia. There is no need to distinguish between the characteristics. (a) Focused on getting 175 in current round with no consideration of longer run implications/impact on other teams' choices in subsequent rounds. (b) Focusing on total earnings in terms of getting 75 each round. No consideration of tradeoffs from choosing C versus D (considering the future). (c) Short sighted—when deciding to defect no consideration of tradeoffs/possible negative consequences of choosing D.

A second coding exercise was conducted to identify strategies from the team dialogues at the *supergame* level. As the name implies, the unit of observation for this exercise was an entire supergame. In extending the game-by-game coding to the supergame level, the simplest case is when a team explicitly stated the strategy underlying their choices. These coders were told to "use the stage game coding as a guide" in determining what strategy to code but were not to be bound by it.[17] If they felt the stage game coding had misidentified the strategy, they should code what they felt was the correct strategy. The instructions to coders stressed that once a strategy was identified, "They don't have to keep saying [their strategy] *as long as their choices correspond to the strategy they had been using.*" Coding for a team only changed if they explicitly stated a new strategy or their actions deviated from their previously stated strategy. When a new strategy was *not* explicitly stated at the time of the change, the coders were instructed to take a holistic approach to determine what strategy was being used. This included looking at the discussions surrounding the point where a change took place along with the team's choices. As noted previously, strategies were often developed over a number of stage games and the clearest statement of a team's strategy was often prior to the time the team started using it. In short, coders used a combination of all the available evidence along with their best judgment to determine the strategy employed.

If a team's strategy changed midway through a long supergame, the supergame was coded based on the *initial* strategy for that supergame. For example, if a team started a supergame with AD but decided in St15 to start playing Lenient Grim, the supergame would be coded AD. If no further discussion of strategies occurred and play in subsequent supergames was consistent with Lenient Grim, these supergames were coded as Lenient Grim.

This extension of the coding to the supergame level was done by two economic graduate student RAs. The coders were given detailed examples of each strategy as well as instructions on how to conduct the coding. They coded the strategies independently and were asked to meet (without the researchers present) to reconcile any differences in their coding. Agreement between the coders was high before they met ($k = 0.79$).[18] The combination of what teams said and did usually made it clear what strategy a team was using.[19] The coders were given the option of leaving a supergame uncoded if they could not identify, or agree, on the strategy, which happened occasionally even after reconciliation.[20]

When choosing the list of strategies to be coded, we drew heavily on the relevant literature in game theory and experimental economics. A reviewer suggested that this may have unwittingly biased the results. To check this, we employed two teams of undergraduate RAs to develop their own lists of strategies with minimal direction.

---

[17] The coding instructions stated, "You will find that occasionally the previous coders got it wrong. So in general anchor off their coding; but if it's obviously off, feel free to change the coding."

[18] The most common disagreements were cases where one coder coded a supergame while the other left it uncoded, and where the coders agreed that some version of Grim was being played but disagreed on which version.

[19] For example, the combination of what teams said and did made AD easy to recognize. In all 183 cases where a team was coded for AD, the team chose D in the first stage game.

[20] There were cases where the two coders did not successfully reconcile their coding. In all such cases, one of the two coders picked a strategy while the other left the supergame uncoded. The analysis that follows uses the strategy from the coder who coded that supergame.

TABLE 3. STRATEGY FREQUENCIES FROM TEAM CHATS

| SG | Always defect | Grim trigger | Lenient grim | Grim w/ counting | TFT | STFT | Complex STFT | Uncoded |
|---|---|---|---|---|---|---|---|---|
| 1 (58 observations) | 55.17% | 22.41% | 3.45% | 3.45% | 6.90% | 1.72% | 6.90% | 0.00% |
| 2 (58 observations) | 53.45% | 12.07% | 8.62% | 8.62% | 8.62% | 3.45% | 3.45% | 1.72% |
| 3 (58 observations) | 48.28% | 22.41% | 10.34% | 1.72% | 3.45% | 5.17% | 6.90% | 1.72% |
| 4 (58 observations) | 43.10% | 24.14% | 20.69% | 3.45% | 1.72% | 1.72% | 3.45% | 1.72% |
| 5 (58 observations) | 36.21% | 29.31% | 12.07% | 8.62% | 5.17% | 5.17% | 3.45% | 0.00% |
| 6 (58 observations) | 34.48% | 25.86% | 18.97% | 5.17% | 5.17% | 6.90% | 3.45% | 0.00% |
| 7 (46 observations) | 17.39% | 36.96% | 17.39% | 4.35% | 8.70% | 4.35% | 6.52% | 4.35% |
| 8 (38 observations) | 7.89% | 39.47% | 23.68% | 7.89% | 10.53% | 5.26% | 5.26% | 0.00% |
| 9 (38 observations) | 7.89% | 42.11% | 18.42% | 13.16% | 7.89% | 5.26% | 2.63% | 2.63% |
| 10 (30 observations) | 16.67% | 30.00% | 20.00% | 6.67% | 6.67% | 13.33% | 3.33% | 3.33% |
| 11 (20 observations) | 20.00% | 25.00% | 20.00% | 5.00% | 5.00% | 20.00% | 5.00% | 0.00% |
| 12 (20 observations) | 15.00% | 25.00% | 25.00% | 5.00% | 10.00% | 10.00% | 10.00% | 0.00% |
| Total (540 observations) | 33.89% | 27.04% | 15.19% | 5.93% | 6.30% | 5.56% | 4.81% | 1.30% |

Their lists were short and included strategies comparable to AD and Grim. They tended to ignore subtleties that are important in understanding behavior; for example, both teams bundled Grim, Grim with Leniency, and Grim with Counting into a single strategy. Details are reported in online Appendix E.

## C. *Teams' Strategies*

Table 3 shows the distribution of strategies identified from the team chats at the supergame level. The number of observations in Table 3 declines in later supergames due to shorter sessions ending. Frequencies are not reported for Always Cooperate or Win-Stay-Lose-Shift as these strategies were never observed.

Initially, Always Defect (AD) was easily the modal choice, averaging 52 percent across the first three supergames. Variants of Grim Trigger (Grim Trigger, Lenient Grim, and Grim with Counting) were also common, totaling 31 percent of the observations in SGs 1–3. The weight on AD decreased continuously across supergames, shifting primarily to variants of Grim Trigger. By SGs 5–7, the last point before substantial dropouts due to sessions ending, the frequency of AD was 32 percent, while variants of Grim Trigger combined for 51 percent of observations.[21] Although never common, Grim with Counting was always present. Variants of tit-for-tat (TFT, STFT, and Complex STFT) were less common than variants of Grim Trigger and did not change frequency with experience (16 percent in SGs 1–3 and SGs 5–7). The movement away from AD toward cooperative strategies parallels the increased cooperation rates in the choice data.

Because the coding in Table 3 was done at the team level, rather than the population level, it is possible to identify how strategies for individual teams changed with

---

[21] For the one session with only six supergames, data from SGs 4–6 were used.

TABLE 4. NUMBER OF SWITCHES BETWEEN CATEGORIES OF STRATEGIES FOR
INDIVIDUAL TEAMS

| Switches | SG1: Always defect | SG1: Cooperative | Row total |
|---|---|---|---|
| 0 | 10 | 16 | 26 |
| | 31.3% | 61.5% | 44.8% |
| 1 | 16 | 2 | 18 |
| | 50.0% | 7.7% | 31.0% |
| 2 | 2 | 6 | 8 |
| | 6.3% | 23.1% | 13.8% |
| 3 | 4 | 1 | 5 |
| | 12.5% | 3.9% | 8.6% |
| 4 | 0 | 1 | 1 |
| | 0.0% | 3.9% | 1.7% |

experience. We classify strategies into two broad categories: AD versus potentially cooperative strategies (variants of Grim Trigger, TFT, or STFT).[22] Table 4 reports how frequently individual teams switched between these two categories across all supergames, broken down by which category the team was in for the first supergame (SG1).

Ten teams played AD for all supergames, and 16 always played one of the cooperative strategies (see the 0-switch row of Table 4). Half of the teams that started out playing AD (50 percent) switched to a cooperative strategy and never switched back (the 1-switch row). Collectively, a clear majority of teams (76 percent) had zero or one switch. Movement was generally away from AD toward cooperative strategies; 63 percent of the teams (20/32) coded for AD in SG1 used a cooperative strategy in their final supergame, but only 12 percent of the teams (3/26) that started with a cooperative strategy in SG1 chose AD in their final supergame. Changes in strategy were largely a one-way street from AD to potentially cooperative strategies.

OBSERVATION 3: *At the supergame level, AD was initially the most frequently used strategy for teams. With experience, cooperative strategies, primarily variants of Grim Trigger, gained weight at the expense of AD. There wasn't much back and forth between AD and cooperative strategies. For the most part, teams either picked AD or one of the potentially cooperative strategies to begin with and stuck with it, or there was one-way movement from AD to a cooperative strategy.*

Table 5 compares strategies identified from team dialogues with the distribution of strategies estimated by the Strategy Frequency Estimation Method (SFEM, Dal Bó, and Fréchette 2011). SFEM is, by far, the most commonly used technique for inferring the strategies underlying choices in IRPD games.[23] SFEM models individuals as playing finite automata, capturing common strategies such as Grim Trigger or Tit-for-Tat. Using a prespecified set of strategies, the model calcu-

---

[22] In the few cases with missing codes, the coding from the previous supergame was employed. There were no cases with more than one consecutive uncoded supergame.

[23] See Dal Bó and Fréchette (2018) for a summary of existing papers that use SFEM.

TABLE 5—STRATEGY FREQUENCIES: CHAT CODING VERSUS SFEM

| Time period | Method | Always defect | Grim trigger | Lenient grim | Counting | TFT | Complex TFT/STFT | STFT |
|---|---|---|---|---|---|---|---|---|
| SGs | Coding | 52.3% | 19.0% | 7.5% | 4.6% | 6.3% | 3.4% | 5.7% |
| 1–3 | SFEM | 43.7% | 2.7% | 0.0% | 0.0% | 28.7% | 5.0% | 19.9% |
| SGs | Coding | 31.6% | 28.7% | 16.7% | 5.7% | 5.7% | 5.7% | 4.6% |
| 5–7 | SFEM | 24.8% | 9.4% | 0.0% | 0.0% | 30.5% | 17.4% | 17.9% |

*Notes:* Lenient Grim combines Grim2 and Grim3. Complex TFT combines TF2T, TF3T, and 2TF2T from SFEM and Complex STFT from the coding. For the session with only six supergames, data from SG4–SG6 are used. Data is included from 348 supergames (coding) and 4,278 choices (SFEM).

lates the likelihood of each agent's observed actions subject to some probability distribution over strategies. The weights on strategies are then fit by maximum likelihood estimation. SFEM is a mixture model; it estimates the probability distribution of strategies across the entire population rather than assigning specific strategies to specific agents. This implies that it cannot identify when individual agents have changed strategies, making an exercise like Table 4 impossible. To ease identification, SFEM is generally estimated on a block of supergames rather than a single supergame. The estimates in Table 5 follow this approach, estimating the model separately for the early (SGs 1–3) and late (SGs 5–7) supergames.[24] Online Appendix F provides an extended discussion of how we fit SFEM.

Both approaches, SFEM and coding from the chats, capture the broad movement away from AD to more cooperative strategies, but there are substantial differences between the distribution of strategies identified by the two approaches. Categorizing contingent strategies as variants of Grim Trigger (Grim Trigger, Lenient Grim, and Grim with Counting) or variants of TFT (TFT, STFT, and Complex TFT/STFT), coding based on the chats always puts less weight on variants of TFT and more on variants of Grim Trigger than SFEM. This is already pronounced in early supergames (SG 1 – 3), where the coding from team chats assigns 16 percent of the population to variants of TFT versus 31 percent to variants of Grim Trigger, compared to 54 percent and 3 percent for SFEM. The difference between methods becomes even larger in late supergames (SG 5 – 7): The coding assigns 16 percent and 51 percent to variants of TFT and Grim Trigger, respectively, as opposed to 66 percent and 9 percent for SFEM.

To understand the broad differences between SFEM and the coding, consider the types of histories that make it possible to distinguish between TFT and Grim Trigger. A sequence like the following must be observed to identify TFT: C/C, C/D, D/D, D/C. There needs to be initial cooperation that falls apart, followed by an attempt to re-establish cooperation. This rather complex set of events occurred infrequently, giving teams little reason to discuss this option before it happened. Teams generally used incomplete strategies, improvising when faced with an unanticipated outcome. SFEM identifies teams' strategies solely based on what

[24] Data from SGs 4–6 were used for the one team session with only six supergames.

TABLE 6—SFEM ESTIMATES, INDIVIDUALS VERSUS TEAMS

| | Individual | | Team | |
|---|---|---|---|---|
| | SGs 1–3 | SGs 5–7 | SGs 1–3 | SGs 5–7 |
| AD | 29.78% | 19.30% | 43.67% | 24.76% |
| AC | 1.11% | 5.85% | 0.00% | 0.00% |
| Grim | 12.16% | 0.00% | 2.69% | 9.37% |
| Lenient grim | 4.88% | 7.81% | 0.00% | 0.00% |
| TFT | 23.32% | 36.98% | 28.74% | 30.53% |
| STFT | 20.66% | 22.10% | 19.89% | 17.94% |
| Complex TFT | 7.05% | 7.96% | 5.01% | 17.40% |
| WSLS | 1.05% | 0.00% | 0.00% | 0.00% |
| p(error) | 7.50% | 5.52% | 3.97% | 2.87% |
| Number of observations | 3,624 | 4,178 | 1,938 | 2,340 |

actions they take, with no means of distinguishing a premeditated choice from improvisation.

OBSERVATION 4: *Compared to SFEM, the coding from team chats identified a relatively higher proportion of variants of Grim Trigger and relatively fewer variants of TFT.*

Comparing strategies between individuals and teams must rely on SFEM rather than the coding. Table 6 compares the distributions of strategies estimated for the individual and team data, subdivided between early (SGs 1–3) and late (SGs 5–7) supergames. To simplify the table, we have combined some strategies: "lenient grim" includes both Grim2 and Grim3, and "complex TFT" includes TF2T, TF3T, and 2TF2T. We also report the estimated probability of an error, defined as playing C when the strategy calls for D or vice versa. Table B1 in online Appendix F includes estimates for the component strategies in these categories, as well as standard errors and the noise parameter.

The estimated distribution of strategies changes over time for both individuals ($\chi^2 = 201.18$; d.f. $= 12$; $p < 0.001$) and teams ($\chi^2 = 230.36$; d.f. $= 12$; $p < 0.001$).[25] In both cases there was a shift away from AD towards more cooperative strategies, but which cooperative strategies gained weight differed. For individuals, almost the entire gain came from variants of TFT (TFT, STFT, and Complex TFT).[26] These strategies gained 16 percentage points, while variants of Grim lost 9 percentage points. Gains were far more even for teams, with variants of TFT gaining 12 percent and variants of Grim growing by 7 percent.

Focusing on the late supergames (SGs 5–7), when agents have had a chance to learn and the strategies employed have settled down, the estimated parameters

[25] Test statistics are log-likelihood ratio tests.
[26] The fitted weight on AC went from 1 percent to 6 percent for individuals. This strategy is ever detected for teams, either by SFEM or the coding.

are significantly different between individuals and teams ($\chi^2 = 35.07$; d.f. $= 12$; $p < 0.001$). This is not due to a dramatic difference in the usage of some particular type of strategy. Instead, it reflects an accumulation of small differences across strategies as well as differing noise parameters.

Given the coding results, it seems likely that SFEM overestimates the frequency of variants of TFT for indindividuals just as these are overestimated for teams. It is therefore difficult to put too much weight on comparing the relative frequencies of particular types of strategies. More telling, the SFEM estimates provide yet more evidence that play by individuals was inherently less stable than team play. The higher frequency of switching reported in Table 2 is paralleled by higher estimated error rates for individuals in SFEM (the error rate gives the probability on playing C when the strategy calls for D, or vice versa). In both early and late supergames, the estimated probability of an error is almost twice as high for individuals versus teams.

OBSERVATION 5: *Based on SFEM, the strategies used by individuals and teams are similar but the error rate is higher for individuals than teams.*

## D. *Coding the Rationale for Switching to Cooperation*

Given that switching to cooperation generally involved extended discussions between teammates, the game-by-game coding isn't terribly useful for understanding *why* teams switched from AD to a potentially cooperative strategy. A separate coding at the *team* level was developed to examine teams' rationale for *unilaterally* switching to cooperation. Two cases were considered: (1) switching from an initial choice of D to an initial choice of C *between* supergames and (2) switching to C following mutual defection *within* the same supergame. There were 42 switches to cooperation between supergames and 24 within supergames, with 39 of the 58 teams having at least one switch.

The authors read through a sample of dialogues when teams switched to cooperation and identified common rationales for the change.[27] The point was not to identify what strategy was used, but rather *why* a switch took place. Two graduate student RAs then coded all of the dialogues in which a switch to cooperation occurred.[28] The instructions the RAs received stressed the need to look not just at the stage game when a change took place, but also the surrounding dialogues, given the extended nature of team discussions. The coders each went through the dialogues separately, and then met to reconcile their codings. The discussion below is based on their reconciled codings. Online Appendix D (Table D2) reports descriptions and frequencies of these coding categories along with Cohen's kappa for agreement rates between the two coders prior to reconciliation. For the most part, the coders agreed reasonably well on the coding.

---

[27] The sample was drawn from the first two team sessions. We used two sessions rather than one to get a reasonably large sample for generating categories.

[28] One of these RAs also did the supergame level coding; due to availability, the other coder was new.

## E. *Understanding the Rationale for Changes in Teams' Choices*

Define a "substantive" discussion as a dialogue for a single stage game with teams discussing what choices to make either in the current stage game or the future (15 percent of stage games). This eliminates dialogues that were unrelated to the experiment. By far the most common topic was what action to take for the current stage game (84 percent of substantive discussions). This occurred in 13 percent of the stage games, as there was little to discuss given that mutual cooperation and mutual defection were quite stable and the interface made it easy to coordinate choices. When there was a genuine need to discuss what action to take, teammates typically did so. In their initial interaction (SG1, St1), 97 percent of teams discussed what action to take. Discussions of what action to use were also common when changing initial actions from the previous supergame (81 percent), changing actions within a supergame (71 percent), or responding to a change in their opponent's action (43 percent).

Discussions of strategies (as opposed to actions) were rarer, occurring in only 21 percent of substantive discussions. Once teams adopted a strategy, they felt little need to continue discussing their strategy if they weren't changing it. Discussions of strategies were surprisingly rare at key moments. A clear majority of teams discussed strategies before the first stage game of the first supergame (69 percent), but relatively few did when switching initial actions at the *start* of a new supergame (34 percent) or when making a unilateral deviation from mutual cooperation or mutual defection *within* a supergame (25 percent). This speaks to the point made previously, that discussions about changing strategies were often extended affairs. When implementing a change in strategy, teammates frequently relied on earlier discussions rather than a discussion at the point in time when the change actually occurred. It follows that we need to use a team's extended conversations when coding their strategies rather than taking a more granular approach.

Teams' discussions provide insight into the motivation behind their initial choices. In SG1, 55 percent of teams were coded as using AD. The most common reason given for this in St1 (42 percent) was not trusting the other team (e.g., "D. You know they'll pick D, so let's get 75 instead of 5."). This was more than twice as frequent as either of the next two most common reasons for choosing AD—myopia (focusing on the short-run benefits of D without recognizing any possible longer-term repercussions) at 20 percent and discussing the impact of current choices on future play at 19 percent.[29] For teams that chose cooperative strategies (variants of Grim or TFT) in SG1, the most common reason given in St1 was discussing the impact of current choices on future play (23 percent) followed closely by discussing the mutual benefits of cooperation (19 percent) and distrust (19 percent). The relationship between what teams discussed and choices in St1 was quite strong as can be seen in the probit

---

[29] It may seem strange that teams choose D while discussing the future negative repercussions of this choice. A team could recognize the negative effect of initially choosing D on future cooperation and still feel that the benefits of protecting themselves against the sucker payoff justified choosing D.

regression below; the dependent variable is a dummy for choosing C in St1 of SG1 (with $p$-values in parentheses).[30]

$$C \ = \ \underset{(0.006)}{-0.753} \, Distrust - \underset{(0.006)}{0.420} \, Myopia + \underset{(0.002)}{0.816} \, MutualBenefits$$

$$- \underset{(0.778)}{0.029} \, DiscussFuture + \underset{(0.331)}{0.135} \, Confusion.$$

OBSERVATION 6: *Distrust of the other team, resulting in fear of getting the "sucker" payoff, was the most common reason coded for justifying an initial choice of AD.*

Switching to a cooperative strategy first involved realizing that the team could do better through mutual cooperation. While this may seem obvious, it was an "aha" moment for many teams, who then had to figure out *how* to coordinate on mutual cooperation.[31] The most common approach was an attempt to lead by example and/or to signal an intent to cooperate ("Leading"). Leading was coded for 53 percent of teams when unilaterally switching to cooperation, being especially common (75 percent) when teams switched to cooperation within a supergame. By cooperating, they hoped their opponent would view them as willing to cooperate and be willing to follow their example. The dialogues are full of examples like the following: "i wonder if choosing C once will make them willing to switch" "I guess we could try?" "it might be worth it" "here we go lol." Teams often viewed leading by example as a way of sending a message in an environment where direct communication was not possible, as one team noted: "this is all so hard without communication" "I know if we could just send them like one sentence we'd have it made."[32]

Leading is a good strategy (in the nontechnical sense) for IRPD games with relatively high continuation probabilities. Because there are, in expectation, many stage games within a supergame, it costs little to take the sucker payoff for a few stage games, compared to the high potential payoff from mutual cooperation in later stage games. Indeed, some teams explicitly discussed the tradeoff between the losses from switching to C and potential gains from achieving mutual cooperation:

> "What if we start with C and do it the whole time" "then we would get 5 points when they choose D" "i think they will start with D, but they are good people so will switch to C" "I don't think they will" "we would make

---

[30] $p$-values are based on robust standard errors. Confusion captures teams that discussed the rules of the game, largely because one of the teammates was confused.

[31] Prior to switching, 34 percent of teams were coded for explicitly discussing the benefits of mutual cooperation and 79 percent of teams were coded as explicitly discussing the need to adopt a new strategy.

[32] There is a related literature on leading by example in public goods literature (for a summary see Cooper and Hamman 2021). The mechanism that makes leading by example successful in public goods game is reciprocity; leaders contribute to the public good, anticipating that the other group members are conditional cooperators and will reciprocate by contributing themselves even though this is *not* an equilibrium strategy (Gächter et al. 2012; Arbak and Villeval 2013). This differs from teams' discussions of leading by example in IRPD games, which were typically framed in terms of sending a message to coordinate on mutual cooperation. Teams were solving a problem of equilibrium selection rather than relying on the kindness of others.

> up the first [stage game] loss quickly" "alright let's do C this round but if
> they choose D i think that's where they'll stay."

Leading was often associated with the lenient grim strategy, as teams understood they needed to give their opponent a chance to catch on. For example, the team just quoted had their opponent choose D in St1 after which the teammate who wanted to cooperate said, "we have to do it twice to see if they change." Their teammate reluctantly agreed, and they achieved mutual cooperation in the next stage game.

A second common rationale for trying cooperation ("learning") was to determine whether the opposing team was willing to cooperate, rather than trying to influence them per se. This was coded for 24 percent of teams unilaterally switching to cooperation, less than half as frequent as leading. The following brief exchange illustrates learning:

> "we can risk it to see what kinda team they are and press C or go safe and
> press D again" "press C" "kk."

Leading and learning were similar rationales for switching to cooperation, primarily distinguished by whether the intent was to influence the other team (leading) or to determine whether the other team were willing to cooperate (learning). It was uncommon to combine the two with only 8 percent of teams coded for both leading and learning.

Only 8 percent of teams explicitly discussed prior play as a reason for switching to cooperation. This is surprising given that teams were more likely to start cooperating if they had recently experienced an opponent who initially cooperated, or following longer supergames with long stretches of mutual defection.[33] The history of prior play presumably prompted teams to think about switching strategies but was not explicitly discussed before switching.

OBSERVATION 7: *The most common approach teams took when trying to coordinate on mutual cooperation was leading by example, signaling their willingness to cooperate.*

Finally, the team dialogues provide clues as to why team choices were so stable compared to individuals' choices. Consider cases where the status quo was either mutual cooperation or mutual defection in the previous stage game. If a team only discussed switching from the status quo, they went through with the switch 54 percent of the time. But when they discussed both switching from the status quo *and* the status quo, they switched only 13 percent of the time. Inertia favored the status quo, consistent with "pluralistic ignorance" noted in the psychology literature (Prentice and Miller 1996). This holds that even when a member of a group privately rejects an opinion or practice, they tend to believe that other members of the

---

[33] Two-thirds of teams that switched to cooperating in St1 experienced an opponent cooperating in St1 of the previous supergame, compared to only one-third two supergames ago. The average length of the supergame preceding the switch was 17.5 stage games, as opposed to 10.8 for two supergames before the switch.

group accept it, making it much easier to abide by an established convention than to change it.

## VI. Finite versus Infinitely Repeated Prisoner's Dilemma Games

This section compares the data here to results from Kagel and McGee (KM; 2016), who ran a series of finitely repeated prisoner's dilemma (FRPD) games using the same stage game payoffs, the same subject population, and the same procedures as those employed here. The number of stage games, ten, matched the expected number of stage games in the IRPD games.[34]

The results for the FRPD games parallel those reported for IRPD games, as mutual cooperation was initially lower for teams than individuals (12.0 percent versus 46.2 percent in the first stage game of SG1) but became greater with experience (52.0 percent versus 38.5 percent in St1 of the final common supergame). As with the IRPD games, mutual cooperation increased faster with experience for teams than individuals ($n = 10$; $z = 1.786$; $p = 0.074$).

Although one might expect less mutual cooperation in St1 for FRPD games than IRPD games, mutual cooperation was *more* common in St1 of SG1 for the FRPD games for individuals (46.2 percent versus 19.2 percent) and teams (12.0 percent versus 10.3 percent). This difference was initially significant for individuals ($n = 11$; $z = 2.64$, $p < 0.01$), but not significant by the final common supergame ($n = 11$; $z = 0.367$; $p = 0.714$).[35] For teams, the difference in St1 cooperation rates between FRPD and IRPD games was not initially significant ($n = 11$; $z = 0.299$; $p = 0.765$). Although the gap grew with experience, it never became statistically significant.[36]

For individuals, mutual cooperation in St1 decreased from 41.0 percent to 31.7 percent between early (SGs 1–3) and late FRPD supergames but increased from 23.1 percent to 29.4 percent for IRPD supergames. Dal Bó (2005) reports similar reductions in St1 cooperation rates for individuals in FRPD games, albeit in much shorter games, commenting that "the effect of the shadow of the future increases with experience." However, this did *not* carry over to teams, as mutual cooperation in St1 increased substantially with experience for *both* FRPD games (from 26.7 percent to 51.0 percent) and IRPD games (from 19.5 percent to 43.3 percent).[37] Teams rapidly learned to cooperate in St1 for FRPD games even though the incentives for starting with C were worse for teams than individuals in early supergames; cooperating in St1 increased expected supergame payoffs by 92 ECUs for individuals versus 22 ECUs for teams.

---

[34] The dataset included five individual sessions and five team sessions. All sessions had at least seven supergames.

[35] To compare apples with apples, these comparisons are based on the first seven supergames in all cases. Dal Bó (2005) reports lower St1 cooperation rates with experience in FRPD games (for individuals) than in parallel IRPD games, but this likely reflects differences in the number of stage games: ten in KM versus two or four in Dal Bó. Embry et al. (2018) report that St1 cooperation rates (for individuals) are increasing in the length of FRPD games.

[36] The difference between FRPD and IRPD games was at its maximum in SG6 (56 percent versus 31 percent), but was still not statistically significant ($n = 11$; $z = 0.739$; $p = 0.460$).

[37] This increase was significant for teams in both FRPD ($n = 5$; $z = 1.761$; $p = 0.078$) and IRPD ($n = 6$; $z = 2.201$; $p = 0.028$) games.

Like the IRPD games, play was more stable for teams than individuals in the FRPD games. The difference in the number of switches looks small in the raw data (1.45 for individuals versus 1.37 for teams), but this is an artifact of differences in the initial conditions. Regressions reported in Table C2 that control for differences in initial conditions between individuals and teams find significantly fewer switches with teams (est. $= -0.165$; SE $= 0.085$; $p = 0.051$). In short, the main differences between teams and individuals in the IRPD games reported here were also present in the FRPD games reported in KM, suggesting that the pattern of results is relatively broad.

Team chats in FRPD games were coded at the stage game level by two graduate students, using essentially the same procedures described above for IRPD games. Like the results reported above for IRPD games, teams that chose D in St1 of SG1 did so primarily out of safety considerations (91.7 percent; 22 of 24 teams choosing D). Like the IRPD games, 70.6 percent (12 out of 17 teams choosing C) did so to elicit mutual cooperation with its higher payoff. Lenient Grim was rarely practiced, likely due to the finite number of stage games. In several cases teams managed to generate cooperation by choosing D in St1, switching to C in St2 if the other team chose C, and continuing with C in St3, hoping to achieve cooperation, similar to the generalized STFT strategy described above.

## VII. Silent Partner Treatment

The comparison of teams and individuals raises a question as to whether behavior differs due to the presence of a teammate per se, possibly due to being responsible for their teammate's earnings, or because of communication and joint decision making between teammates. To distinguish between these two hypotheses, a silent partners treatment was implemented. Like the team treatment, subjects were assigned to fixed two-subject teams at the beginning of the experiment, with payoffs the same for both teammates as a consequence of their choices. One member of the team was randomly chosen for the role of Decider and the other was assigned the role of Silent Partner. These roles were fixed for the duration of the experiment. All decisions for the team were made by the Decider. The Decider knew they were responsible for their Silent Partner's payoffs. The Silent Partner observed the Decider's choices and outcomes, but there was no communication between the two.[38]

Four silent partner sessions were conducted. The same procedures were employed as in the main sessions, and the same seeds were used as in the final four team sessions. All sessions ran for at least as long as the matched team session, with the data reported on below restricted to common supergames. Redoing the regressions reported in online Appendix C with only data from the four seeds used in the silent partner treatment, the differences between individuals and teams reported in Observations 1 and 2 remain statistically significant.

---

[38] To make it less obvious which subjects were Deciders, silent partners were encouraged to make choices indicating what "they would have [done] if they were the Decider." These played no role in determining the team's choice and Deciders could *not* observe their silent partner's choices. Most silent partners made choices (88 percent of all stage games), but the correlation with the Deciders' choices was not especially high ($r = 0.47$, subject to making a choice).

The silent partner treatment was designed to test how much of the difference between teams and individuals was due to having a partner per se, rather than joint decision making and communication. This might matter because decision makers might become more risk averse if their decision affects the risks born by a passive second party (as in Bolton, Ockenfels, and Stauf 2015), making them less willing to initially take the risk of cooperating. Additionally, if Deciders exert greater cognitive effort on choosing a strategy because of other-regarding preferences (their effort now benefits their Silent Partner as well as themselves), they should learn to cooperate faster in the silent partners treatment than in the individual treatment. The greater stability of play by teams grows from the process of joint decision making – when the teammates disagree, there is a strong tendency to stick with the status quo. Merely having a partner should not have the same effect.

Limiting the dataset to the four common seeds, initial mutual cooperation rates and the changes in cooperation rates over time differed little between individuals and the silent partner treatment. In the first stage game of SG1, the mutual cooperation rate for the silent partners treatment was 25 percent, compared with 26 percent for individuals and 16 percent for teams. Turning to the final common supergame, mutual cooperation in the first stage game was 56 percent for the silent partners treatment, an increase of 31 percentage points. The analogous mutual cooperation rates were 44 percent and 68 percent for individuals and teams, representing increases of 18 and 53 percentage points. The direction of these effects for the silent partner treatment relative to individuals were the same as in the team treatment, but the magnitudes were far smaller.

Table C1 reports probit regressions estimating the size of these effects, controlling for the seed, supergame, length of the previous supergame, and experience with cooperation by opponents. None of the effects of the silent partners treatment, relative to individuals, were statistically significant.

The average number of switches per supergame was 1.89 for the silent partners treatment. This is higher than the average number of switches in the individual treatment (1.29) rather than lower, as in the team treatment (0.87). That said, the estimated difference between the silent partners treatment and the individual treatment is not significant (see Table C2). Once again, play in the silent partner treatment differed little from play by individuals.

The similarities between the silent partners treatment and the individual treatment are consistent with the idea that the differences between individual and team decision making were primarily due to joint decision making and the associated communication between teammates. Having a partner per se could potentially have an effect through either risk or social preferences, but this was not the case.

The effects of joint decision making and team communication are inherently difficult to disentangle since it is not possible to coordinate choices without some type of communication between teammates. To the extent that this has been studied, full bilateral communication is a necessary condition for the strong performance of teams relative to individuals. Neither unilateral communication nor communication limited to an exchange of proposed actions with no further explanation is sufficient to replicate the effect of team play with free-form communication (Cooper and Kagel 2016; Arad, Grubiak, and Penczynski 2021).

## VIII. Discussion and Conclusions

There were two motivations for the experiment reported here: (1) To compare the behavior of individuals and freely interacting two-person teams in IRPD games with perfect monitoring; and (2) To use the dialogues between teammates to understand teams' underlying behavioral processes. Arguably, the greatest value of studying teams is use of their conversations to understand how and why decisions were made. Economists have become increasingly interested in process data (e.g., fMRI, eye tracking, reaction times) to understand decision making, including analysis of team chat. Team discussions are a natural part of the decision making process and provide direct insights into how and why decisions come about.

Teams were less cooperative than individuals in early supergames, but cooperation rates increased more rapidly for teams, resulting in significantly more cooperation than individuals in late supergames. Additionally, team play was more stable both *within* supergames and between supergames. This implies that both mutual cooperation and mutual defection were more persistent within supergames for teams.

Analysis of teams' dialogues provided direct insights into the strategies teams used and the rationale underlying their choices. Discussions that explicitly laid out strategies in the game theoretic sense were rare. Rather than specifying plans for all possible histories, strategies typically developed over time and were improvisational in nature. Further, it was not uncommon for teams to switch strategies midway through a supergame. Teams' choices were not arbitrary, since the dialogues showed that they generally thought carefully about how to play the game, but their thought processes were far less structured than modelers typically assume. Given the broad similarity of team and individual behavior, it seems reasonable to assume that similar processes underlie individual choices.[39]

Related to the latter point, results from Romero and Rosokha (2019) are consistent with individuals improvising within supergames. Their experiment used individuals as agents, long IRPD games ($\delta = 0.98$), and direct elicitation of strategies. There were frequent changes within supergames (2.37 per supergame) when subjects could costlessly change their strategy within a supergame, consistent with the improvisation observed in the team chats.[40]

Strategies identified from the team dialogues show that teams moved steadily across supergames from always defecting (AD) to cooperative strategies with little backsliding. Comparing strategies identified from the team dialogues with the distribution of strategies estimated by SFEM, coding assigns more weight to variants of Grim and less to variants of TFT. Rather than deciding in advance how to handle unlikely contingencies, teams adopted simpler strategies like Grim and its variants and then adjusted on the fly to unexpected circumstances.

---

[39] Play by individuals and teams clearly is not identical. That said, it is qualitatively similar with cooperation rates starting low and growing with experience. Likewise, we argue that the **underlying processes are likely to be qualitatively similar as well.** For example, individuals, like teams, no doubt rely largely on improvisation as opposed to well thought out state contingent plans.

[40] Improvisation is somewhat different than simply changing strategies midway through a supergame. Teams in our experiments generally did not have fully specified strategies, and filled in the details in response to the actions of their opponents. Because Rosokha and Romero used direct elicitation, agents always had a complete strategy.

The most common reason teams gave for choosing to defect in early supergames was distrust of their opponent. The most frequent approach used when switching to a cooperative strategy was leading by example—cooperating for a few initial stage games to signal their willingness to cooperate, hoping that their opponent would catch on. Another common approach was learning about their opponent, where teams initially cooperated in an attempt to find out if their opponent was willing to cooperate. These two approaches, leading and learning, are distinguished by the motivation for cooperating, either trying to influence their opponent's behavior or testing for their responsiveness to cooperation.

The analysis of team dialogues has important implications for how cooperation emerges and suggests new ways of modeling this process. There have been relatively few learning models used to study IRPD games (see Dal Bó and Fréchette 2011 and Romero and Rosokha 2019 for exceptions), but the strong dynamics in the data point to the importance of such models. Teams that led by example were consciously trying to affect their opponent's choices.[41] This suggests that learning models for IRPD games need to incorporate strategic sophistication since trying to anticipate and influence the decision process of other players is the essence of strategic sophistication.[42] Teams often improvised, exploring new strategies. Models that treat teams as choosing a strategy at the beginning of each supergame and sticking with it miss much of what teams actually did. They observed what happened, thought about it, possibly engaged in some experimentation, and often adjusted their strategy. A good model of learning in IRPD games would capture this continuous process of experimentation and learning.

Teams started out cooperating less than individuals consistent with the "discontinuity effect" reported in the psychology literature for PD games. That literature attributes the lower cooperation rates for teams to greater distrust and support for self-serving choices. The data here are consistent with this explanation, although bounded rationality (myopia) also played an important role in initial decisions to defect. What the psychology literature fails to identify, due to typically implementing a single supergame in an experimental session, is that with experience the cooperation rate for teams surpasses that of individuals.

Finding that teams' choices were more stable than individuals' choices was not expected. Two explanations for this suggest themselves. First, this might reflect teams being more rational than individuals. In line with this, Proto, Rustichini, and Sofianos (2019), who investigated repeated games where individuals were stratified into higher and lower cognitive ability cohorts, reported that higher cognitive ability individuals were more stable in their choices (specifically, less likely to deviate following mutual cooperation in an IRPD game). A second possibility is that team choices were more stable because team decision making was inherently biased in favor of the status quo. The dialogues support this as teams that discussed both switching strategies and the status quo generally stuck with the status quo. Switching choices was difficult for

---

[41] This relates to strategic teaching as documented by Hyndman, Terracol, and Vaksmann (2009) in coordination games and Hyndman et al. (2012) in normal form games with a unique Nash equilibrium which is on the Pareto frontier.

[42] See the SEWA model of Camerer, Ho, and Chong (2002) for an example of how sophistication can be added to a model of learning.

teams because unanimity was necessary to make a change, consistent with the pluralistic ignorance literature in psychology (Prentice and Miller, 1996).

In previous work (Cooper and Kagel, 2005) comparing strategic play by individuals and teams, we have used the truth wins (TW) model of Lorge and Solomon (1955). This was developed to model decision making by teams for logic problems that have a demonstrably correct solution, with the reference point being that a team solves the logic problem if any of its members, working independently, would have solved the problem. The TW model does *not* apply to IRPD games as there is no singular optimal strategy (i.e., a demonstrably correct solution). Instead, the optimal strategy varies depending on the strategies adopted by other players. However, playing C in St1 leads to higher average payoffs than playing D (see fn 10). Based on this, define play of C in St1 as the "empirically optimal" action. Under a version of the TW model using this alternative definition of optimality, team play is initially below the TW benchmark as teams initially cooperate less than individuals. Over time, as teams become more cooperative, team play just catches up with the lower range of the 90 percent confidence interval for this TW benchmark. The relatively fast growth of cooperation in the team data is not captured by the TW model, consistent with our observation that the TW model is not well suited for IRPD games. Our interpretation of the data is that the comparison between teams and individuals reflects a combination of two factors. Initially, what dominates is the discontinuity effect (i.e., teams are more sensitive to the risk of the other side defecting, and are therefore less willing to cooperate than individuals). However, teams are faster than individuals to learn how to build cooperation through tactics like leading by example. The TW model is designed to capture the adoption of eureka insights, such as realizing that leading by example is a good approach, but is poorly suited for factors like matters of judgment (e.g, risk preferences) or the lack of an unambiguously optimal strategy in IRPD games. See online Appendix G for a fuller discussion of the TW model.

The results reported here are based on a single specific environment, with relatively long supergames and perfect monitoring. Both of these features likely affected the findings. The popularity of leading by example was no doubt due in part to the length of the supergames. With a continuation probability of 0.90, it was cheap to send the opposing team a message by choosing C for one or two stage games, compared to the potential benefits of inducing them to cooperate. Likewise, the lack of noise made it less important for teams to think about what to do under various contingencies. This may change with imperfect monitoring as teams respond to the inherent uncertainty in their opponents' actions. The point will not be just to see how teams' choices change as the structure of the game varies, but to use their discussions to understand the process underlying their choices.

## REFERENCES

**Aoyagi, Masaki, V. Bhaskar, and Guillaume R. Fréchette.** 2019. "The Impact of Monitoring in Infinitely Repeated Games: Perfect, Public, and Private." *American Economic Journal: Microeconomics* 11 (1): 1–43.

**Arad, Ayala, Kevin P. Grubiak, and Stefan P. Penczynski.** 2022. "Does Communicating within a Team Influence Individuals' Reasoning and Decisions?" *Experimental Econonomics* https://doi.org/10.1007/s10683-022-09786-3.

**Arbak, Emrah, and Marie Claire Villeval.** 2013. "Voluntary Leadership: Motivation and Influence." *Social Choice and Welfare* 40: 635–62.

**Aumann, Robert J., and Lloyd S. Shapley.** 1994. "Long-Term Competition – A Game Theoretic Analysis." In *Essays in Game Theory in Honor of Michel Maschler*, edited by Nimrod Megiddo, 1–15. New York, NY: Springer.

**Blonski, Matthias, Peter Ockenfels, and Giancarlo Spagnolo.** 2011. "Equilibrium Selection in the Repeated Prisoner's Dilemma: Axiomatic Approach and Experimental Evidence." *American Economic Journal: Microeconomics* 3 (3): 164–92.

**Breitmoser, Yves.** 2015. "Cooperation, but No Reciprocity: Individual Strategies in the Repeated Prisoner's Dilemma." *American Economic Review* 105 (9): 2882–2910.

**Camerer, Colin F., Teck-Hua Ho, and Juin-Kuan Chong.** 2002. "Sophisticated Experience-Weighted Attraction Learning and Strategic Teaching in Repeated Games." *Journal of Economic Theory* 104 (1): 137–88.

**Casari, Marco, Jingjing Zhang, and Christine Jackson.** 2016. "Same Process, Different Outcomes: Group Performance in an Acquiring a Company Experiment." *Experimental Economics* 19 (4): 764–91.

**Cason, Timothy N., and Vai-Lam Mui.** 2019. "Individual versus Group Choices of Repeated Game Strategies: A Strategy Method Approach." *Games and Economic Behavior* 114: 128–45.

**Cooper, David J., and John R. Hamman.** 2021. "Leadership and Delegation of Authority." In *The Handbook of Labor, Human Resources, and Population Economics*, edited by Marie Claire Villeval and Klaus Zimmerman. Berlin: Springer-Verlag.

**Cooper, David J., and John H. Kagel.** 2005. "Are Two Heads Better Than One? Team versus Individual Play in Signaling Games." *American Economic Review* 95 (3): 477–509.

**Cooper, David J., and John H. Kagel.** 2016. "A Failure to Communicate: An Experimental Investigation of the Effects of Advice on Strategic Play." *European Economic Review* 82: 24–45.

**Cooper, David J., and John H. Kagel.** 2023 "Replication data for: Using Team Discussions to Understand Behavior in Indefinitely Repeated Prisoner's Dilemma Games." American Economic Association [publisher], Inter-university Consortium for Political and Social Research [distributor]. https://doi.org/10.38886/E176881V1.

**Cooper, David J., Ian Krajbich, and Charles N. Noussair.** 2019. "Choice-Process Data in Experimental Economics." *Journal of the Economic Science Association* 5: 1–13.

**Dal Bó, Pedro.** 2005. "Cooperation under the Shadow of the Future: Experimental Evidence from Infinitely Repeated Games." *American Economic Review* 95 (5): 1591–1604.

**Dal Bó, Pedro, and Guillaume R. Fréchette.** 2011. "The Evolution of Cooperation in Infinitely Repeated Games: Experimental Evidence." *American Economic Review* 101 (1): 411–29.

**Dal Bó, Pedro, and Guillaume R. Fréchette.** 2018. "On the Determinants of Cooperation in Infinitely Repeated Games: A Survey." *Journal of Economic Literature* 56 (1): 60–114.

**Dal Bó, Pedro, and Guillaume R. Fréchette.** 2019. "Strategy Choice in the Infinitely Repeated Prisoners' Dilemma." *American Economic Review* 109 (11): 3929–52.

**Davis, James H.** 1992. "Some Compelling Intuitions about Group Consensus Decisions, Theoretical and Empirical Research, and Interpersonal Aggregation Phenomena: Selected Examples 1950–1990." *Organizational Behavior and Human Decision Processes* 52 (1): 3–38.

**Embrey, Matthew, Guillaume R. Fréchette, and Sevgi Yuksel.** 2018. "Cooperation in the Finitely Repeated Prisoner's Dilemma." *Quarterly Journal of Economics* 133 (1): 509–51.

**Fehr, Ernst, and Paul M. Glimcher.** 2013. *Neuroeconomics: Decision Making and the Brain*. London, UK: Academic Press.

**Feri, Francesco, Bernd Irlenbusch, and Matthias Sutter.** 2010. "Efficiency Gains from Team-Based Coordination—Large-Scale Experimental Evidence." *American Economic Review* 100 (4): 1892–1912.

**Fischbacher, Urs.** 2007. "z-Tree: Zurich Toolbox for Ready-Made Economic Experiments." *Experimental Economics* 10 (2): 171–78.

**Friedman, James W.** 1971. "A Non-cooperative Equilibrium for Supergames." *Review of Economic Studies* 38 (1): 1–12.

**Fudenberg, Drew, and Eric Maskin.** 1986. "The Folk Theorem in Repeated Games with Discounting or with Incomplete Information." *Econometrica* 54 (3): 533–56.

**Gächter, Simon, Daniele Nosenzo, Elke Renner, and Martin Sefton.** 2012. "Who Makes a Good Leader? Cooperativeness, Optimism, and Leading-By-Example." *Economic Inquiry* 50 (4): 953–67.

**Greiner, Ben.** 2015. "Subject Pool Recruitment Procedures: Organizing Experiments with ORSEE." *Journal of the Economic Science Association* 1: 114–25.

**Harsanyi, John C., and Reinhard Selten.** 1988. *A General Theory of Equilibrium Selection in Games*. Cambridge, UK: MIT Press.

**Hyndman, Kyle, Antoine Terracol, and Jonathan Vaksmann.** 2009. "Learning and Sophistication in Coordination Games." *Experimental Economics* 12 (4): 450–72.

**Hyndman, Kyle, Erkut Y. Ozbay, Andrew Schotter, and Wolf Ze'ev Ehrblatt.** 2012 "Convergence: An Experimental Study of Teaching and Learning in Repeated Games." Journal of the European Economic Association 10 (3): 573–604.

**Kagel, John H., and Peter McGee.** 2016. "Team versus Individual Play in Finitely Repeated Prisoner Dilemma Games." American Economic Journal: Microeconomics 8 (2): 253–76.

**Kocher, Martin G., and Matthias Sutter.** 2005. "The Decision Maker Matters: Individual versus Group Behavior in Experimental Beauty-Contest Games." Economic Journal 115 (500): 200–23.

**Lorge, Irving, and Herbert Solomon.** 1955. "Two Models of Group Behavior in the Solution of Eureka-Type Problems." *Psychometrika* 20 (2): 139–48.

**Lugovskyy, Volodymyr, Daniela Puzzello, and James Walker.** 2018. "On Cooperation in Finitely and Indefinitely Repeated Prisoner's Dilemma Games." Unpublished.

**Maciejovsky, Boris, Matthias Sutter, David V. Budescu, and Patrick Bernau.** 2013. "Teams Make You Smarter: How Exposure to Teams Improves Individual Decisions in Probability and Reasoning Tasks." *Management Science* 59 (6): 1255–70.

**Prentice, Debra A., and Dale T. Miller.** 1996. "Pluralistic Ignorance and the Perpetuation of Social Norms by Unwitting Actors." *Advances in Experimental Social Psychology* 28: 161–209.

**Proto, Eugenio, Aldo Rustichini, and Andis Sofianos.** 2019. "Intelligence, Personality and Gains from Cooperation." *Journal of Political Economy* 127 (3): 1351–90.

**Rand, David G., Joshua D. Greene, and Martin A. Nowak.** 2012. "Spontaneous Giving and Calculated Greed." *Nature* 489: 427–30.

**Romero, Julian, and Yaroslav Rosokha.** 2018. "Constructing Strategies in the Indefinitely Repeated Prisoner's Dilemma Game." *European Economic Review* 104: 185–219.

**Romero, Julian, and Yaroslav Rosokha.** 2019. "The Evolution of Cooperation: The Role of Costly Strategy Adjustments." *American Economic Journal: Microeconomics* 11 (1): 299–328.

**Wildschut, Tim, and Chester A. Insko.** 2007. "Explanations of Interindividual-Intergroup Discontinuity: A Review of the Evidence." *European Review of Social Psychology* 18 (1): 175–211.

**Wildschut, Tim, Brad Pinter, Jack L. Vevea, Chester A. Insko, and John Schopler.** 2003. "Beyond the Group Mind: A Quantitative Review of the Interindividual-Intergroup Discontinuity Effect." *Psychological Bulletin* 129 (5): 698–722.