

# Implicit and Explicit Communication in Experiments

A DISSERTATION  
SUBMITTED TO THE FACULTY OF THE GRADUATE SCHOOL  
OF THE UNIVERSITY OF MINNESOTA  
BY

**Piotr Evdokimov**

IN PARTIAL FULFILLMENT OF THE REQUIREMENTS  
FOR THE DEGREE OF  
DOCTOR OF PHILOSOPHY

Advisers: David Rahman and Aldo Rustichini

June 2014

© Piotr Evdokimov 2014

## Acknowledgments

I am immeasurably thankful to my advisers for showing me what it means to be a good economist. To Aldo, a special word of thanks for teaching me how to run an experiment and demonstrating by example that rigorous thinking and openness to new ideas are not mutually exclusive. David, I cannot imagine having a better friend, teacher, and colleague than you. I am also grateful to Itai Sher, Jan Werner, and Kim Sau Chung for a myriad of constructive comments; to my fellow students for invaluable hours spent discussing economics (Jan Grobovsek, Zachary Mahone, Thomas Youle, and Jan Duras, this means you guys); and to my family, without which I wouldn't be here.

## Dedication

Dear esteemed, so cherished everyone! No disguised, onerous line in this thesis, locked ever closely, lastingly, obdurately up, dared appear; nor drudging hurdles overturn, veritable excesses roll back, exogenously. Fellowships offer relief; elementary tokens help enormous encumbrances yield ever slightly. Obrigado, friends that helped! Eternal Love. And to V. N., for teaching me what an acrostic is.

## Abstract

This thesis is composed of three chapters, distinct in focus but unified by the basic question of how communication affects strategic behavior. The first chapter deals with forward induction, a common equilibrium selection principle in game theory which relies on a sort of *implicit* communication. The second studies communication in a face-to-face environment. The third chapter explores the effects of minimal, explicit communication in a repeated game with frequent actions and imperfect monitoring.

Several general conclusions can be drawn from the work presented below. The first is that the sort of implicit, rational communication that economic theory describes does not seem to be widespread. Only about 30% of the subjects who participated in the experiments described in Chapter 1 demonstrated strong evidence of understanding the forward induction principle. The second is that face-to-face communication can have adverse effects on an individual's earnings that are not predicted by economic theory. Thus, it is shown in Chapter 2 that agreeable workers are paid less by their managers in a controlled bargaining environment. The third is that minimal, explicit communication which in theory has no bearing on behavior improves the earnings of individuals a repeated interaction with imperfect monitoring. On the other hand, information management institutions that in theory can allow players to sustain more cooperation can in practice have adverse effects.

The subject of Chapter 1 is implicit communication in a strategic environment. Imagine that two players, Person 1 and Person 2, are engaged in a strategic interaction with two stages, and that Person 1 made a move in the first stage. Forward induction is the notion that from this first stage move Person 2 can infer what Person 1 believes will happen later. Theoretical models of forward induction take a prominent place in the literature (Battigalli and Siniscalchi, 2002; Kohlberg and Mertens, 1986), but whether or not people make forward induction-like inferences in the laboratory is an open question. In the experiments described in Chapter 1, detailed reports from participants playing a battle of the sexes game with an outside option. Approximately a third of these reports exhibited an excellent understanding of forward induction, and these reports were associated more strongly with forward induction-like behavior than reports consistent with first mover advantage and other reasoning processes. The experiments also provide some evidence that forward induction thinking can be learned through observation of other players' actions.

In most laboratory experiments, subjects interact anonymously through a computerized interface. Communication in the field, however, often takes place in environments where agents interact face-to-face. In Chapter 2, I report the results of a study that allowed this kind of interaction to test the hypothesis that an individual's economic choices and outcomes are significantly affected by the personality of others, measured in terms of the "Big Five" (Costa and McCrae, 1992). Specifically, the study introduced a bargaining experiment where output was produced by a worker, a manager determined how this output was divided between the worker and herself, and the two parties interacted face-to-face to determine their bargaining positions. The results of this study attribute a significant effect of the worker's personality on her bargaining power.

Chapter 3 studies the effects of canonical information management institutions on cooperation in a game with imperfect monitoring. Delay of information is the first institutional manipulation this study considers; the second is the ability of players to communicate their strategies, and the third is bounded rationality in the form of constraints on reaction time. The results of this study show that subjects earn significantly more *without* delay of information, a result that cannot be explained by standard repeated games models, that communication always improves welfare, and that average payoffs in one of our treatments (with communication and no delay) are significantly greater than the upper bound on public Nash equilibrium payoffs. Exploring the possibility that this is driven by bounded rationality in the form of reaction lags, the study finds that slowing down the experiment has no significant effect on behavior.

# Contents

<b>Acknowledgements</b>	<b>i</b>
<b>Dedication</b>	<b>ii</b>
<b>Abstract</b>	<b>iii</b>
<b>List of Tables</b>	<b>viii</b>
<b>List of Figures</b>	<b>ix</b>
<b>1 Forward Induction</b>	<b>1</b>
1.1 Introduction . . . . .	1
1.1.1 Literature Review . . . . .	4
1.2 Experimental Design . . . . .	6
1.2.1 Protocol . . . . .	6
1.2.2 Feedback Structure . . . . .	7
1.2.3 Questions . . . . .	8
1.2.4 Treatments . . . . .	10
1.3 Rating Procedure . . . . .	11
1.4 Results . . . . .	14
1.4.1 Behavior in the Game . . . . .	14
1.4.2 Report Evaluations . . . . .	18
1.4.3 Relationship Between Thinking and Behavior . . . . .	21
1.5 Discussion and Robustness . . . . .	23

1.6	Conclusion . . . . .	24
<b>2</b>	<b>Personality and Bargaining Power</b>	<b>26</b>
2.1	Introduction . . . . .	26
2.2	Literature on Income and Personality . . . . .	29
2.3	Experimental Design . . . . .	31
2.3.1	Details of the Experimental Design . . . . .	33
2.4	Results . . . . .	35
2.4.1	Income, Inequality and Incentives . . . . .	35
2.4.2	Correlations With Personality . . . . .	38
2.4.3	Identification . . . . .	39
2.4.4	Evaluation of Others Through Questionnaires . . . . .	43
2.4.5	Robustness Checks . . . . .	45
2.4.6	The Worker's Effort . . . . .	45
2.5	Conclusion . . . . .	47
<b>3</b>	<b>Cooperative Institutions</b>	<b>49</b>
3.1	Introduction . . . . .	49
3.2	Related work . . . . .	53
3.2.1	Experimental literature . . . . .	53
3.2.2	Theoretical literature . . . . .	54
3.3	Experimental design . . . . .	55
3.4	Theoretical predictions . . . . .	59
3.4.1	Public equilibrium payoffs . . . . .	59



3.4.2	How information delay can help . . . . .	60
3.4.3	Delay with practical cut-offs . . . . .	62
3.4.4	How bounded rationality can help . . . . .	63
3.5	Results . . . . .	64
3.5.1	Dynamics . . . . .	68
3.5.2	Strategies . . . . .	70
3.5.3	Periodicity of behavior . . . . .	74
3.6	Discussion . . . . .	76
	<b>References</b>	<b>79</b>
<b>4</b>	<b>Appendix</b>	<b>90</b>
4.1	Questions for Second Movers . . . . .	90
4.2	Instructions . . . . .	92
4.3	Quiz questions for participants . . . . .	96
4.4	Summary of the Data . . . . .	96
4.5	Additional Robustness Checks . . . . .	99
4.6	Instructions . . . . .	101
4.7	Personality Questionnaire . . . . .	110
4.8	Proofs . . . . .	113
4.9	Instructions to treatment NC . . . . .	115

# List of Tables

1	Table 1 . . . . .	15
2	Table 2 . . . . .	17
3	Table 3 . . . . .	18
4	Table 4 . . . . .	20
5	Table 5 . . . . .	22
6	Table 6 . . . . .	36
7	Table 7 . . . . .	37
8	Table 8 . . . . .	39
9	Table 9 . . . . .	40
10	Table 10 . . . . .	44
11	Table 11 . . . . .	46
12	Table 12 . . . . .	47
13	Table 13 . . . . .	65
14	Table 14 . . . . .	66
15	Table 15 . . . . .	69
16	Table 16 . . . . .	71
17	Table 17 . . . . .	73
18	Table 18 . . . . .	98
19	Table 19 . . . . .	100
20	Table 20 . . . . .	100

# List of Figures

1	Figure 1 . . . . .	2
2	Figure 2 . . . . .	16
3	Figure 3 . . . . .	19
4	Figure 4 . . . . .	42
5	Figure 5 . . . . .	58
6	Figure 6 . . . . .	60
7	Figure 7 . . . . .	61
8	Figure 8 . . . . .	67
9	Figure 9 . . . . .	72
10	Figure 10 . . . . .	75
11	Figure 11 . . . . .	76
12	Figure 12 . . . . .	97
13	Figure 13 . . . . .	98

# 1 Forward Induction<sup>1</sup>

## 1.1 Introduction

Forward induction (FI) is the notion that what an individual does in an early stage of a multi-stage interaction contains information about what he or she will do later. This information can be used strategically by others to facilitate coordination. FI was introduced to economics by [Kohlberg and Mertens \(1986\)](#) and found several successful applications, particularly in models of industrial organization ([Ponssard, 1991](#); [Ben-Porath and Dekel, 1992](#); [Bagwell and Ramey, 1996](#)). In the [Bagwell and Ramey \(1996\)](#) model, two firms, an entrant and an incumbent, make sequential investment capacity decisions in the presence of multiple equilibria, with the entrant moving first. The incumbent, assuming that the entrant is rational, deduces the entrant's post-entry production from the latter's capacity investment. If the capacity commitment is such that the entrant is not able to recoup its investment with a market sharing equilibrium, the incumbent concludes that the entrant will produce at natural monopoly levels and shuts down.

The experimental evidence of FI has been mixed. The [Bagwell and Ramey \(1996\)](#) capacity model has been tested in the laboratory ([Brandts et al., 2007](#)); other laboratory studies of FI include [Cooper et al. \(1993\)](#), [Brandts and Holt \(1995\)](#) and [Huck and Müller \(2005\)](#). All of these studies show that while subjects often coordinate on equilibria selected by FI, other factors, such as asymmetries between the players, play a substantial role in determining behavior. In a game where forward and backward induction make different predictions, [Balkenborg \(1994\)](#) showed that the FI outcome is selected less than 20% of the time. Summarizing the state of the literature, [Samuelson \(2005\)](#) reported that “the experimental evidence has not been particularly supportive of forward induction, suggesting that theories based on forward induction could well be reconsidered.” This suggestion notwithstanding, developments in the theoretical literature on FI have anything but slowed: [Govindan and Wilson \(2009\)](#), [Battigalli and Friedenbergl \(2012\)](#), [Müller \(2012\)](#), and [Man \(2012\)](#) are some recent examples. It seems that the jury on FI is still out.

We contribute to the experimental literature on FI by incorporating a new source of data: subjects' elicited beliefs, as expressed in reports of a relatively free-form

---

<sup>1</sup>Joint work with Aldo Rustichini.

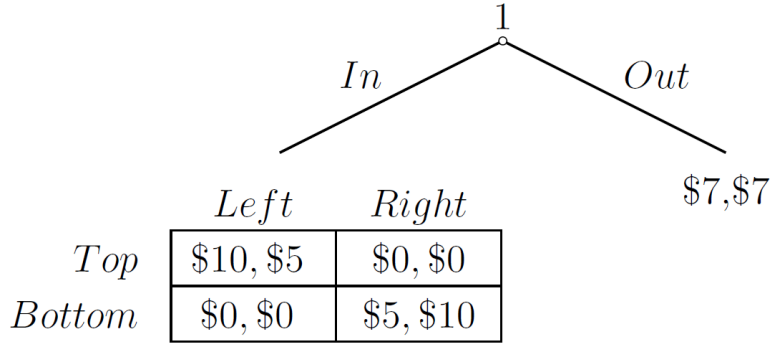


Figure 1: **The game used in the experiment.**

nature. All previous experimental studies FI have focused on participants’ choices, leaving unsettled the question of whether FI thinking is actually used. Our experiment is the first to address the issue directly. We designed it around a battle of the sexes game with an outside option, the simplest game in which FI makes a prediction (Figure 1). The first mover chooses between an outside option of \$7 and playing a battle of the sexes game, in which the first mover’s favored payoff is \$10. Under these parameters, the first mover’s decision to play in the subgame contains the following implicit message, quoted from Kohlberg and Mertens (1986): “Look, I had the opportunity to get [\$7] for sure, and nevertheless I decided to play in this subgame, and my move is already made. And we both know that you can no longer talk to me, because we are in the game, and my move is made. So think now well, and make your decision.” For the second mover to “think well” is to realize that any strategy other than a sure choice of *Top* is inconsistent with a rational first mover’s decision to play in the subgame. The best response of a second mover who “thinks well,” therefore, is to choose *Left* whenever he or she given the opportunity to play. Accounting for the above, a rational first mover chooses *In* in the first stage of the game, and *Top* in the second.

We instructed participants to play this game in anonymous, randomly determined matches. In several rounds, participants answered a series of questions aimed at assessing one’s understanding of the FI argument. Their answers to these questions were scored by independent evaluators following an incentivized procedure described in Xiao and Houser (2005), described in more detail in Section 1.3. The evaluators found a substantial fraction of reports, namely 30%, to show an excellent understanding of FI, and 33% to show only a *possible* understanding thereof.<sup>2</sup> Reports in the

<sup>2</sup>These “weaker” reports were consistent with FI as well as, for example, the first mover advantage

“stronger” category were associated with more FI-like choices in the game. This was true in rounds when reports were provided, as well as those when they were not. In the former case, a “stronger” report increased the probability of FI-like behavior by 36%, while a “weaker” one increased this probability by 19%. These effect sizes were significantly different ( $P < 0.01$ ). Further, participants who understood FI exhibited a significantly greater *change* in behavior—in the direction of putting a higher probability on  $((In, Top), Left)$ —than that experienced by their peers. Thus, a significant portion of forward induction-like behavior was driven by forward induction reasoning.

Our results provide some evidence that FI reasoning can be learned. While reports in late periods of the experiment did not exhibit greater understanding of FI, the more one’s partners chose to play in the subgame, the more likely the player was to produce a FI report. We find, however, that observation of other players needs a substantial period of time for a participant’s thinking about the game to change.

One of our treatments allows us to distinguish between not understanding FI but lacking confidence that others possess such understanding. Thus, a subject might say that the second mover will choose *Right* if the subgame is played, but that they themselves would choose *Left* in the role of the second mover, since the first mover’s decision to forgo the outside option contains an implicit message. We identify a small proportion of participants—approximately 9%—whose reasoning follows such a framework. These are participants that understand FI, but do not think that others are able or willing to do the same.

Finally, we find that participants’ chosen actions were broadly consistent with previous studies in the literature. If the subgame was played, the FI action was selected in a large majority of cases, 95% of the time for first and 82% of the time for second movers. Behavior exhibited substantial learning ( $P < 0.001$ ), and in rounds where participants reported on their thinking, FI outcomes were significantly more likely ( $P < 0.001$ ). This raises the question of how general our conclusions about the prevalence of FI thinking are. We stress that although we cannot say how often such thinking is used in absence of belief elicitation, we can put a lower bound on how many participants (approximately 30%) have the capacity to understand FI, and study the implications of such understanding for behavior. This is the primary purpose of our work. The effect of belief elicitation in our experiment is broadly in line with existing studies, which are reviewed below.

---

hypothesis.

### 1.1.1 Literature Review

Kohlberg and Mertens (1986) associated FI with stability of equilibria. Van Damme (1989) suggested that FI should instead correspond to the following requirement, which is not always associated with stability: “...in generic two-person games in which player  $i$  chooses between an outside option or to play a game  $G$  of which a unique viable equilibrium  $e^*$  yields this player more than the outside option, only the outcome in which  $i$  chooses  $G$  and  $e^*$  is played in  $G$  is plausible.” Cho and Kreps (1987) showed stability to be useful for studying signaling games. Other theoretical papers that follow the Kohlberg and Mertens (1986) approach include De Sinopoli (2004), Glazer and Weiss (1990), Ponssard (1991), Ben-Porath and Dekel (1992), Bagwell and Ramey (1996), Gul and Pearce (1996), Hauk and Hurkens (2002) and Govindan and Wilson (2009).

More recently, Battigalli and Siniscalchi (2002) associated FI with the epistemic notion of Strong Belief in Rationality (SBR). Stated informally, “player  $i$  *strongly believes* that an event  $E \neq \emptyset$  is true (i.e., adopts  $E$  as a ‘working hypothesis’) if and only if she is certain of  $E$  at all histories consistent with  $E$ ” (Battigalli and Siniscalchi, 2002, p. 357). To see how SBR drives FI reasoning, consider initial (first stage) beliefs of players in the game illustrated by Figure 1. If the first mover expects the second to choose *Right* in the subgame, then *Out* is a rational action. Hence, even a second mover with strong belief in rationality can initially believe that the first mover thinks that the second will choose *Right*. If the first mover chooses *In*, however, this belief is no longer viable: SBR forces the second mover to now take the view that their opponent assigns a high enough probability to the action *Left* that choosing *In* is justified. Thus, SBR and observed first stage behavior jointly determine beliefs that the second mover has in the subgame. Recent papers that treat FI within the epistemic framework developed by Battigalli and Siniscalchi (2002) include Chung (2011), Battigalli and Friedenberg (2012), and Müller (2012).

Cooper et al. (1993) reported the results of the first laboratory experiment on FI, focusing on behavior in a battle of the sexes game with an outside option. Other experimental studies of FI include Van Huyck et al. (1993), Balkenborg (1994), Brandts and Holt (1995), Cachon and Camerer (1996), Muller and Sadanand (2003), Huck and Müller (2005), Brandts et al. (2007), Blume and Gneezy (2010), and Shahriar (2013). Most of these studies suggest that FI has predictive power, but identify alternative behavioral forces that are able to account for observed behaviors. Several

papers, including [Cooper et al. \(1993\)](#), point to asymmetry between the first and the second mover as a partial explanation of their results. One basic conceptual hurdle for choice-based studies is the following. Even if FI is shown to have an effect in a symmetric game, it's not clear what fraction of behavior is explained by FI in a game that's asymmetric. Additional data are needed in order to tease apart FI from alternative behavioral explanations. Such is the approach we pursue.

There is a growing experimental literature on belief elicitation, the basic findings of which are reviewed in [Schotter and Trevino \(2014\)](#). Most papers in this literature use quadratic scoring rules; we do the same when eliciting point beliefs. Although there is some evidence that these rules are associated with misreporting ([Armantier and Treich, 2013](#)), there is also evidence that the distortions in reports are small or non-existent ([Sonnemans and Offerman, 2001](#)). Some studies, in contexts other than FI, have addressed the question of how participants' reported thinking relates to behavior. Most of these studies (using quadratic scoring rules) have shown that beliefs to play a significant role in determining subjects' actions. This is the conclusion of [Nyarko and Schotter \(2002\)](#), [Danz et al. \(2012\)](#), and [Hyndman et al. \(2012\)](#), for instance. One prominent exception is [Costa-Gomes and Weizsäcker \(2008\)](#), who conclude that "stated beliefs reveal deeper strategic thinking than [...] actions. On average, [subjects] fail to best respond to their own stated beliefs in almost half of the games." Unlike [Costa-Gomes and Weizsäcker \(2008\)](#), we find that subjects tend to best respond to their stated beliefs and that strategic sophistication is highly predictive on behavior.

Other papers have looked at the effect of eliciting incentivized and non-incentivized beliefs on behavior. While the results are mixed, several studies, such as ([Croson, 2000](#)), have found that incentivized belief elicitation has a significant effect on subjects' choices. Ours also finds this effect to be significant. [Schotter and Trevino \(2014\)](#) interpret the existing evidence to suggest that to the extent belief elicitation alters behavior, it makes stable, best response behavior emerge sooner. Our results are in part consistent with this interpretation: the prevalence of FI-like behavior grows over time even in absence of belief elicitation, but asking participants questions makes such behavior more likely. We also, however, observe *transient* effects of belief elicitation on choice. This suggests that the learning-based explanation of [Schotter and Trevino \(2014\)](#) is incomplete.

Since thought processes other than FI can lead to FI-like beliefs, we collect free-



form written report data in addition to point beliefs of the kind that quadratic scoring rules elicit. These beliefs were scored following the method described in [Xiao and Houser \(2005\)](#). Specifically, we hired outside evaluators, instructed them how to classify the free-form reports, and incentivized performance using a coordination game in which a rater is paid if his or her classification agreed with that of the majority. Most of the experimental studies in economics that utilized free-form report data focused on the effect of communication on behavior ([Charness and Dufwenberg, 2006](#); [Schotter and Sopher, 2007](#); [Xiao and Houser, 2005](#); [Houser and Xiao, 2011](#)); no paper prior to ours used such methods to investigate FI thinking.

## 1.2 Experimental Design

We now outline the experimental design. [Section 1.2.1](#) provides the details of the protocol; [Section 1.2.2](#) focuses on our choice of feedback structure; [Section 1.2.3](#) describes the questions participants answered, and [Section 1.2.4](#) motivates and discusses our treatments. A description of how the reports were categorized is deferred to [Section 1.3](#).

### 1.2.1 Protocol

The experiment included four treatments (labeled T1-T4), all run in the Social and Behavioral Sciences Laboratory at University of Minnesota in the fall semesters of 2011 (T1 and T2) and 2012 (T3 and T4). The treatments were programmed with z-Tree ([Fischbacher, 2007a](#)) and used University of Minnesota undergraduate students as participants. Upon signing a consent form, each subject was seated at a computer terminal and given a paper copy of the instructions, the complete text of which is included in [Appendix 4.9](#). Subjects first read the instructions, then completed a multiple choice quiz on the details of the game, displayed on a computer screen. Completing the quiz correctly on the first try resulted in a small (\$1) reward. Then, forty (in T1 and T2) or eight (in T3 and T4) rounds of the game were played.

Half of the participants in each session played as first movers in even and second movers in odd rounds; for the other half, the reverse was true, and this information was known to all subjects. Hence, every session was split into two groups, and, throughout the experiment, participants that started out as first movers were matched

with those that started out in the other role. Subjects were not told how many times they will play the game, or *when* they will have to report on their thinking, but they were informed that they will “play the game a number of times,” that their partner will be randomly selected in every round, and that they will answer questions about their thinking at some points in the experiment. Some of these answers were incentivized, as described in Section 1.2.3. Second movers were always informed of their opponent’s first stage choice, that is, whether their partner had chosen *In* or *Out*, but no other feedback was given. Participants in T1 and T2 were paid on the basis of the outcome of four rounds, and those in T3 and T4 were paid for the outcome of two. These rounds were drawn at random at the end of the session in addition to a \$5 show up fee. Participants received their earnings from the games at the end of the session, and, four to five weeks later, a check in the mail for their reports.<sup>3</sup>

### 1.2.2 Feedback Structure

As described above, our participants were provided with limited feedback about the game: Second movers only knew if the subgame was played, and first movers had no feedback at all. We consider this a crucial element of the design. [Kahneman \(2011, p. 203\)](#), for example, notes that people “are prone to assess the quality of a decision not by whether the process was sound but by whether its outcome was good and bad.” Because we viewed FI as precisely the kind of sound process described by [Kahneman](#), we needed to ensure that it is not confounded by the observation of the choices of other players. We expected this feature of our design to have consequences on behavior. On the one hand, it is argued by [Rick and Weber \(2010\)](#) that “withholding feedback encourages deeper thinking about the game.” On the other, not knowing what their opponents choose in the subgame, subjects who don’t understand FI, or lack confidence in the ability of others to use FI reasoning, are likely to assume the worst (namely a choice of *Left*) and opt for the outside option. Thus, one would expect more outside options being chosen in an experiment with little opportunity for learning.

---

<sup>3</sup>The time interval was necessary to rate the reports.

### 1.2.3 Questions

We now describe the questions answered by first movers in rounds when reports were provided. For reasons described in [Section 1.3](#), the answers to questions posed to second movers were not scored by the evaluators. These answers did, however, play a role in determining first movers' payments, following scoring rules described below. We include second movers' questions in [Appendix 4.1](#).

Before the decision to go *In* or *Out* was made, each first mover answered F1-F3:

F1. If you go in, what move will the second player make?

Possible answers: *Left* for sure, 90/10 *Left/Right*, 80/20 *Left/Right*, ..., *Right* for sure.

F2. Why will the second mover make this move?

F3. If you go in, what will the second mover think you will do?

Possible answers: *Top* for sure, 90/10 *Top/Bottom*, 80/20 *Top/Bottom*, ..., *Bottom* for sure.

If *Out* was chosen in the first stage, the first mover saw a screen that said: "We are now informing the second mover that you chose *Out*. Press OK to continue." If *In* was chosen, the first mover had to choose their next move and answer F4.

F4. Why did you make this move?

We provided incentives for multiple choice and typed answers with monetary rewards. Where participants guessed actions that were eventually realized, we used the quadratic scoring rule to calculate their payment (denoted by  $\pi$  below). For example, if a first mover chose *In*, they were paid in cents for their answer to F1 according to

$$\pi_{S1} = \begin{cases} 100 - (\text{Prob}_{F1}(\textit{Right}))^2 & \text{if the second mover chose } \textit{Left} \\ 100 - (\text{Prob}_{F1}(\textit{Left}))^2 & \text{if the second mover chose } \textit{Right} \end{cases},$$

where  $Prob_{F1}(Right)$  is the probability assigned to the second mover choosing *Right* by the first mover's answer to F1. This scoring rule ensures that a risk-neutral decision maker reports truthfully.

If a player made a guess about what their opponent is thinking, they were fined with a square of the distance between their and their opponent's reports. For example, if the first mover chose *In*, the payment for F3 was calculated as

$$\pi_{F3} = 100 - (S3_{Top} - F3_{Top})^2,$$

where  $S3_{Top}$  is the probability assigned by the second mover to their partner choosing *Top* in their answer to S3. This scoring rule ensures that a subject with point beliefs reports his or her conjecture, while a subject maximizing expected earnings with a non-trivial belief set reports his or her expectation of the partner's report.<sup>4</sup>

To provide appropriate incentives for the written reports, we told participants that each statement describing *a partner's* thoughts or behavior would be evaluated by us, and classified into a category. The same procedure would be applied to the partner's description of *his own* reasoning process, and only the former subject (the one describing his partner's thoughts) would be paid in the event the categories matched. A statement describing one's own behavior was not incentivized. The specific instructions were:

For each question in which you explain what the other player is doing (or thinking), this player will have a corresponding question in which they have to explain their own behavior (or thoughts). We will evaluate your answer by placing it into one of several categories, and do the same for the answer the other player provides. If the categories match, we will pay you \$1. For example, if your explanation of why the other player will behave in a particular way matches that player's explanation of their behavior, you get \$1. If your explanations do not match, we will pay you \$0. If one of you does not provide a short answer, we will pay you \$0.

Moreover, the subjects were instructed as follows:

---

<sup>4</sup>If  $f(x)$  be a subject's subjective probability distribution over *stated* beliefs of his partner, the minimization problem is:

$$\min_{p \in [0, 1]} \int (x - p)^2 f(x) dx.$$

It is trivial to check that  $p = \int x f(x) dx$  satisfies the first and second order conditions.

We will not pay you additional money for answers to questions about your own thoughts and behavior, but obtaining your considered opinion is important for our study, so please be as detailed as possible in your answers.

#### 1.2.4 Treatments

In treatment **T1**, subjects answered the questions described in [Section 1.2.3](#) in rounds 5, 6, 11, 12, 17, and 18 of the game. In treatment **T2**, the questions were answered in rounds 25, 26, 31, 32, 37, and 38. Thus, each subject in these treatments had three opportunities each to describe their thinking as first and second mover. In addition to addressing the basic questions of how prevalent FI thinking is and how relates to behavior, these two treatments allowed us to see test whether such thinking is learned on a round-by-round basis. A cross-treatment comparison allowed us to see whether participants were more likely to demonstrate FI thinking if they were given twenty rounds of additional experience with the game.

One may fail to understand the logic of FI altogether, or doubt that his or her opponent thinks about the game in such terms. The two treatments described below were designed to determine whether confidence in one's partner played a role in shaping participants' reported beliefs. **T3** and **T4** each consisted of eight rounds of play, with questions answered in the fifth and sixth rounds. Both treatments differed from T1 in two respects: session length and payment structure. Thus, subjects were paid for the outcomes of *two* games in T3 and T4, whereas in T1 they were paid for four. Treatment T4 had the additional feature that first movers choosing *Out* were asked the following three questions:

F5. What would the second mover choose if instead you had moved in?

Possible answers: *Left* for sure, 90/10 *Left/Right*, 80/20 *Left/Right*, ..., *Right* for sure.

F6. Why would the second mover have made this choice?

F7. Put yourself for a moment in the role of your partner. Suppose you were the second mover, and you saw the first mover move in. What would YOU think in that case?

Under the null hypothesis of confidence in one’s opponent not being a factor, the estimate of the number of FI thinkers obtained from T4 should be identical to that obtained from T3. Rejection of the null would suggest that the number of FI thinkers estimated in the other treatments is biased downward.

### 1.3 Rating Procedure

The reports were rated by outside evaluators.<sup>5</sup> The evaluators were recruited in January 2013 with an e-mail to PhD students at the School of Physics at the University of Minnesota. Students of physics were chosen for two reasons. Firstly, they were unacquainted with the researchers, so their reports were unlikely to be biased. Secondly, the quantitative nature of their discipline increased the likelihood that the raters will understand the notion of FI sufficiently well to provide reasonable ratings. For their participation, the raters were paid \$15 per hour. In addition, four reports were selected at random at the end of the session, and each rater was paid an additional \$10 for each instance of their classification agreeing with that of the majority.<sup>6</sup>

Due to the subjective nature of the procedure, it was important that the same raters classified every report in the data set. Moreover, the reports needed to be rated in a single session. Allowing the raters to go home and return at a different date to finish the scoring procedure carries a cost of losing control over what the raters say to each other. Thus, for instance, the raters might have agreed on a scheme that maximizes payment at the expense of classification accuracy. On the other hand, 389 reports were provided by first movers alone in T1-T4.<sup>7</sup> We estimated the rating procedure to take at least three hours even if reports of second movers are not shown to the raters. Thus, for the sake of time and simplicity, only reports of first movers were scored. We take the view that thinking about the game is a unitary process: if a player understands FI as a first mover, he should understand it as a second mover, as well. While the latter sort of understanding requires a smaller degree of strategic

---

<sup>5</sup>In addition, a portion of the reports was scored by the authors of the paper, as in [Charness and Dufwenberg \(2006\)](#). Our ratings agree with those of the physics students to a large extent; [Section 3.6](#) discusses this in more detail.

<sup>6</sup>One rater asked if such a scheme incentivizes reporting what one believes his peers to believe, rather than one’s personal beliefs. Our answer was that our true interest lied in each rater’s *own* beliefs about the reports.

<sup>7</sup>The maximum possible number of reports is  $48 \times 3 + 60 \times 3 + 52 + 70 = 446$  reports, we have fewer because some subjects failed to provide typed answers to our questions.

sophistication, our own preliminary investigations suggest that this is not the case: We find as many subjects showing evidence of FI as first but not second movers as those demonstrating the opposite pattern. This point is elaborated upon in [Section 3.6](#).

After entering the laboratory and signing a consent form, the raters were instructed by a PowerPoint-based presentation.<sup>8</sup> The presentation described the game, the structure of the experiment producing the report data (feedback, order of choices, questions answered by participants, etc.), and explained the notion of FI. It then gave examples of real reports provided by first movers that could be grouped in each of the three suggested categories: those showing evidence of understanding FI (“YES” reports), those *potentially* showing such evidence (“MAYBE” reports), and those *not* showing it (“NO” reports).

A report consisted of the first mover’s answers to F1, F2, and F4 (the latter if available), as well as the player’s choice in the game. The raters were told that a report cannot be judged to show understanding of FI without an explanation of behavior (i.e., an answer to F2 or F4) present, since a belief that the second mover will choose *Right* in the subgame with 100% certainty and, in turn, believe that the first mover will chose *Top* can be supported by thought processes other than FI. Thus, the raters were told to pay attention to the text and only look at the answer to F1 or the subject’s choice if something in a typed response needs to be clarified. Importantly, the raters were also told that a player choosing *Out* can understand FI, and that their job as raters is to identify such understanding—not its reflection in behavior.

A YES report was explained to be one that explicitly stated that giving up an outside option of \$7 implies that the first mover is looking to make \$10 and/or choose *Top*, or strongly hinted at the understanding of this fact. The example given to the raters contained the following sentence (here, as in other examples, original spelling and grammar are preserved):

I chose to go in so I have sent the message to the second mover that I intend to make more than the original \$7, I can only do this if I choose top.

If the described reasoning process was consistent with understanding FI, but could

---

<sup>8</sup>The slides are available upon request.

have been used in a similar game with the outside option replaced by (\$3, \$3), the raters were told that the report belongs in the MAYBE category. Two examples of MAYBE reports were given, one of them replicated below.

F1. If you go in, what move will the second player make?

Answer: 80/20 *Left/Right*.

F2. Why will the second mover make this move?

Answer: I figure that the second mover will make this move because they will figure that I will choose top, because that way I make more money, and instead of not making any money, the second mover will most likely choose left and get \$5.

*(In, Top)*

F4. Why did you make this move?

Answer: I chose top because that way I will make more money, but I am hoping that the second player will choose left because they know that I will choose top only because I can make more money and I hope they think that I am not being generous by choosing bottom for them to make \$10 instead of me.

Finally, the raters were told that a NO report should be one clearly not consistent with understanding of FI. Two examples of such reports were given, one of them stating the probability of the second mover choosing *Right* in the subgame to be 70%, and containing the following passage.

Realistically speaking, I think it's 50/50 in terms of what the smartest move is. Once you figure in the human emotion known as hope, preferences for certain options show up. So the idea behind my 70/30 is that hope can account for 20 percent change. I'm really just going on hunches. As for right, most people are right handed. So I put two and two together.

Following the instructions, the raters were told that they will be presented with the reports. They were also told that after all the reports are rated, the experiment will have a second, shorter, portion, and that the instructions to this half of the experiment will be given after everyone is done with the initial set of ratings. The



reports were then presented through an E-Prime interface. A rater had a maximum of 40 seconds to rate each report. An untimed break was provided every 20 trials. During the break, the rater could rest, go to the bathroom, or, if desired, continue the scoring procedure.

After every rater finished with the first set of reports, a second set of instructions, describing T4, was presented. The raters were instructed that they will score 38 reports consisting of the answers to F5, F6 and F7, and that they should follow the criteria established previously. After the second set of reports was rated, four reports were chosen, and each rater’s classifications were compared to that of the majority. The raters were then paid for their time and classification choices.

## 1.4 Results

We now turn to our experimental findings. In [Section 1.4.1](#), we summarize participants’ behavior and analyze their choices in the game with panel data methods. In [Section 1.4.2](#), we estimate how many participants understood FI and study how likely this kind of thinking is to be learned. [Section 1.4.3](#) addresses the relationship between reported thinking and observed behavior.

### 1.4.1 Behavior in the Game

A total of 230 subjects participated in the experiment. [Table 13](#) shows the breakdown of the overall sample size by treatment, as well as a summary of choice data. The fraction of first movers choosing *In* fluctuated from 33% in T3 to 45% in T1, and the overall average (39%) was substantially smaller than the 80% observed by [Cooper et al. \(1993\)](#). We take this to be a consequence of the structure of feedback in our design. This conclusion is supported by the observation, discussed in [Section 1.2.2](#), that subjects who do not observe their partner’s choice in the subgame (and hence lack evidence of second movers choosing *Left*) are more likely to choose the outside option than those who do.<sup>9</sup> [Table 13](#) also tabulates behaviors of first and second

---

<sup>9</sup>Another difference between our and [Cooper et al.](#)’s design concerns structure of payment. In [Cooper et al. \(1993\)](#), participants played for points, and instead of our \$7, \$10 and \$5, they could earn 300, 600 and 200 points, respectively, in each game played; the probability of receiving a payment was calculated as the number of points earned divided by 1000. We chose dollar amounts for two reasons. First, [Selten et al. \(1999\)](#) showed that lotteries do a poor job at inducing risk neutrality.

movers in the case the subgame was played. Overall, the probability of first movers choosing *Top* was 95%, while that of second movers choosing *Right* was 82%. Second movers in our experiment chose the FI strategy less often than those in Cooper et al. (1993), presumably because exclusion of feedback boosted their optimism about attaining the more desired outcome. Comparing behavior in the first eight rounds across treatments, we found no differences in probabilities of first movers choosing *In* ( $P = 0.553$ ) or *Top* ( $P = 0.291$ ).<sup>10</sup> Although we found a treatment effect on the choices of second movers ( $P < 0.05$ ), this effect lost significance when T2 was excluded from the analysis ( $P = 0.202$ ). We conclude that subjects in different treatments behaved similarly in early periods of the game.

	T1	T2	T3	T4	T1-T4 (all)
N	48	60	52	70	230
In as first mover	0.45	0.36	0.33	0.35	0.39
In as first mover	0.39	0.33			
Top as first mover in subgame	0.98	0.95	0.85	0.88	0.95
Top as first mover in subgame	0.95	0.88			
Left as second mover in subgame	0.86	0.84	0.74	0.65	0.82
Left as second mover in subgame	0.77	0.85			
FI behavior	0.43	0.34	0.25	0.26	0.36
FI behavior	0.34	0.28			

Table 1: **Summary of behavioral data.** For T1 and T2, the table shows the overall likelihood of behavior (topmost values in each panel), as well as the likelihood in the first eight periods (bottom). Thus, second movers in T1 choose *Left* 86% of the time overall, and 77% of the time in early periods of the experiment.

In what follows, we focus on “periods” of play. A period is a pair of two consecutive rounds (e.g., period 1 consists of rounds 1 and 2). Thus, participants in T1 answered questions about their thinking in periods 3, 6, and 9; participants in T2 answered them in periods 13, 16, and 19, and participants in T3 and T4 answered them in period 3. Notice that every participant played once as a first mover and once as a second mover in every period.<sup>11</sup> A participant is defined to be a *FI player* in a given

Second, we wanted to keep our instructions as simple as possible.

<sup>10</sup>Recall that the eighth round was final in T3 and T4.

<sup>11</sup>We focus on periods in subsequent analysis because, as discussed in Section 1.3, we consider

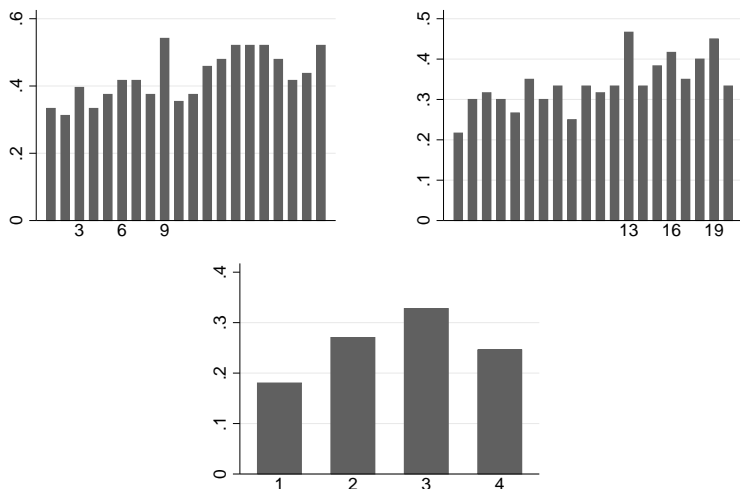


Figure 2: **The probabilities of displaying FI-like behavior in different treatments.** Note the spikes around periods when participants answered questions about their thinking.

period if, in this period, he chose *In* and *Top* as a first mover, and, provided that he was given an opportunity to play in the subgame, *Left* as a second mover. *FI player* should be read as shorthand for “a player exhibiting behavior consistent with FI.” We emphasize that being such a player does not imply an understanding of the FI argument: in [Section 1.4.3](#), we study the relationship between thinking and behavior directly. The bottom rows of [Table 13](#) summarize the prevalence of FI-like behavior in T1-T4. Overall, behavior was consistent with being a FI player 36% time, and no significant treatment effect was found in the first four periods ( $P = 0.159$ ). In [Figure 2](#), we plot the probabilities of being a FI player in each period of the game for participants in T1 (a), T2 (b) and T3 and T4 (c). The figure suggests that this probability went up in periods participants answered questions about their thinking.

To confirm these observations, we estimated the following logit model:

$$P(FIPLAYER_{it} = 1) = \frac{1}{1 + \exp^{-1} \left( \beta_0 + \beta_1 vr_{it} + \beta_2 t + \sum_{k=2}^5 \beta_{3,k} Treatment_{ik} \right)}, \quad (1.1)$$

---

thinking about the game to be a unitary process. See also [Section 3.6](#) for evidence that this assumption is reflected in the report data.

where  $vr_{it} = 1$  if  $i$  provided a verbal report in period  $t$  and 0 otherwise, and

$$FIPLAYER_{it} = \begin{cases} 1 & \text{if, in period } t, i \text{ chose } (In, Top) \text{ as first} \\ & \text{and, given the chance, } Left \text{ as second mover .} \\ 0 & \text{otherwise} \end{cases}$$

Here, as in the rest of our analysis, the standard errors were clustered by subject. The results showed  $\beta_1$  to be positive and highly significant ( $P < 0.001$ ; results reported in Table 2). It is perhaps surprising that FI-like behavior became less likely in each period that followed one with reports.<sup>12</sup> We take this to indicate that our questions encouraged participants to think more carefully about the game and one’s behavior therein, and that their thinking reverted to a baseline, less focused, level when questions were not asked. The coefficient on period number was also significant, suggesting that FI-like behavior is learned over time ( $P < 0.001$ ).

Dependent variable=FI behavior	
Period no.	0.007**** (0.002)
Verbal report in period	0.076**** (0.017)
2.Treatment	-0.091 (0.072)
3.Treatment	-0.131 (0.076)
5.Treatment	-0.119* (0.062)
Observations	2648

Table 2: **Panel analysis of behavioral data.** Marginal effects of logit regressions. Subject-clustered standard errors in parentheses. \*  $P < 0.1$ , \*\*  $P < 0.05$ , \*\*\*  $P < 0.01$ , \*\*\*\*  $P < 0.001$ . Participants exhibit learning of FI behavior, which in addition becomes more likely in periods when verbal reports are provided.

<sup>12</sup>In Figure 2, this can be seen as “spikes” in FI-like behavior. In the above-described regression, this can be inferred from the fact that learning is controlled for.

## 1.4.2 Report Evaluations

Overall, the evaluators judged 30% of reports to show an understanding of FI, placing 33% of reports in the MAYBE and 37% in the NO category. Table 3 shows how these numbers break down across treatments. Recall that all treatments but T2 had subjects report on their thinking in the third period (T1, in addition to this, asked for reports in periods 6 and 9). It is apparent from the table that the fraction of period 3 reports judged to show evidence of FI fluctuated across treatments (from 15% in T3 to 28% in T1 and T4; see Table 3); however, there was no significant treatment effect on how reports in this period were categorized ( $P = 0.255$ ). In T1 and T2, where participants provided reports in multiple periods of the game, 34% of reports were judged to show an understanding of FI, and the fraction of *participants* showing understanding of FI was even higher. Seven subjects provided no reports at all in these two treatments, and, of the remaining 101, 47% had at least one YES report.

	T1	T3	T4	Period 3 (all treatments)	
YES	0.28	0.15	0.28	0.24	
MAYBE	0.38	0.31	0.27	0.31	
NO	0.35	0.54	0.45	0.45	
	Period 3	Period 6	Period 9	T1 (all periods)	
YES	0.28	0.37	0.37	0.34	
MAYBE	0.38	0.32	0.34	0.34	
NO	0.35	0.32	0.29	0.32	
	Period 13	Period 16	Period 19	T2 (all periods)	
YES	0.34	0.33	0.33	0.34	
MAYBE	0.38	0.31	0.37	0.36	
NO	0.28	0.35	0.30	0.31	

Table 3: Percentages of reports in different categories, as scored by outside evaluators.

One explanation of these findings is that a large fraction of participants ( $\approx 50\%$ ) understood FI but wavered in their confidence of other subjects' understanding, resulting in a smaller fraction of YES reports in any given period. The treatment T4 was designed to address this hypothesis; in it, first movers choosing *Out* were asked to explain how *they* would have thought about the game as a second mover. Of the 38

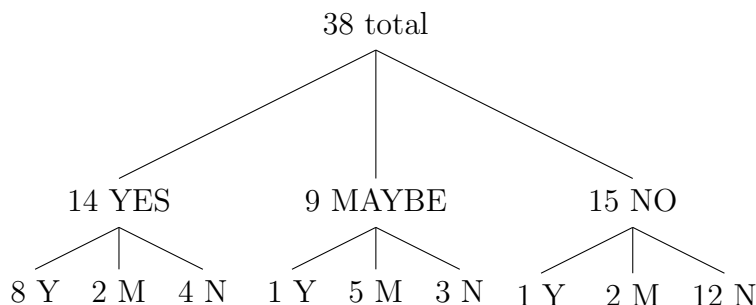


Figure 3: **The 38 auxiliary reports of T4 and the basic reports associated with them.**

reports, 14 were judged to show evidence of FI reasoning (Figure 3). Of these 14, six were provided by participants judged to provide a MAYBE or NO report while they were playing the game. Thus, six out of 70 participants in T4, approximately 9%, understood FI but lacked the confidence to express this understanding while they were making their first stage decisions. These participants, while few, may explain the discrepancy between the percentages of players who understood FI in the longer (T1 and T2) and shorter (T3 and T4) treatments: If we only consider reports obtained in the course of the game, we find that a larger fraction of participants understood FI in T1 and T2 than in T4 ( $P < 0.05$ ). However, if we account for reports of players choosing *Out*, we find that 36% of participants in T4 show understanding of FI, and that this proportion is not significantly different from the 47% estimated in T1 and T2 ( $P = 0.208$ ).

Treatments T1 and T2 allowed us to explore another possibility—that participants in the longer treatments were more likely to understand FI because they were provided with an opportunity to learn it. Table 3, panels b and c, shows how reports were classified in different periods. An exact Fisher test shows no significant effect of period number on report categories in T1 ( $P = 0.897$ ) or T2 ( $P = 0.933$ ). Recall that reports were elicited early on (periods 3, 6, and 9) in T1 and later in T2 (period 13, 16 and 19). If experience with the game made FI thinking more likely, we should observe a significant difference between how reports are categorized in T1 and T2. However, no such difference is manifest in our data ( $P = 0.988$ ).<sup>13</sup>

We thus hypothesized that FI was learned on the basis of observing other players,

<sup>13</sup>In light of the fact that FI-like behavior became more prevalent over time (Table 2), these findings may appear puzzling. We discuss the issue of learning further in Section 1.4.3.

	YES report	MAYBE report	NO report
No. of times partner chose In	0.035*** (0.013)	0.010 (0.014)	-0.046*** (0.016)
Period no.	-0.008 (0.009)	-0.005 (0.010)	0.014 (0.009)
2.Treatment	-0.032 (0.115)	0.047 (0.122)	-0.015 (0.111)
3.Treatment	-0.175** (0.083)	-0.027 (0.089)	0.202** (0.096)
5.Treatment	-0.030 (0.084)	-0.077 (0.078)	0.107 (0.086)
Observations	389	389	389

Table 4: **The effect of feedback on the thinking of participants.** Marginal effects on predicted probabilities in a multinomial logit regression. Subject-clustered standard errors in parentheses. \*  $P < 0.1$ , \*\*  $P < 0.05$ , \*\*\*  $P < 0.01$ , \*\*\*\*  $P < 0.001$ .

and not playing the game per se, and estimated a multinomial logit model with

$$p_{itj} = \frac{\exp(\beta_{0j} + \beta_{1j}t + \beta_{2j}OBS_{it} + \sum_{k=2}^5 \beta_{3jk}Treatment_{ik})}{1 + \sum_{j \in \{YES, MAYBE\}} \exp(\beta_{0j} + \beta_{1j}t + \beta_{2j}OBS_{it} + \sum_{k=2}^5 \beta_{3jk}Treatment_{ik})},$$

where  $p_{itj}$  is the probability of subject  $i$  producing a report of type  $j$  in period  $t$ , and  $OBS_{it}$  is the number of times  $i$  observed his or her partner choose  $In$  in periods prior to  $t$ . The model showed that observing one's partner choose  $In$  made a YES report more and a NO report less likely ( $P < 0.01$  for both marginal effects; results shown in Table 4). However, when we estimated the same model in three separate specifications, one restricting the sample to T1, the other to T2, and the third to T3 and T4, we found these effects with  $P$ -values below 0.001 for T2, but not in the other treatments (largest  $P = 0.346$ ).<sup>14</sup> Thus, a substantial amount of feedback was required for a participant's thinking about the game to change.

<sup>14</sup>Inclusion of interactions of  $OBS_{it}$  with the treatment dummies in the full model produces the same conclusion.

### 1.4.3 Relationship Between Thinking and Behavior

In rounds when reports were provided, a YES report was associated with FI-like behavior 69%, a MAYBE report 50%, and a NO report 16% of the time. Moreover, these proportions are significantly different ( $P < 0.001$  for all pairwise comparisons). To quantify the effect of providing each kind of report on behavior, we added YES, MAYBE, and NO report dummies to the model specified by Equation 1.1, restricting the data to periods in which reports were provided.<sup>15</sup> The results, collected in Table 5, column 1, show that a YES report increased the likelihood of being a FI player by 36%, while a MAYBE report increased it by 19%. Crucially, the effect of a YES report was significantly greater than that of a MAYBE report ( $P < 0.01$ ). A better understanding of FI was associated with more FI behavior.

It was also true that subjects who produced at least one YES report exhibited FI-like behavior in more periods of the game *without reports* than subjects who did not ( $P < 0.001$ ). This was true even when we restricted our attention to participants who did not provide any reports in the NO category ( $P < 0.001$ ). To further investigate this relationship, we controlled for participants' reported thinking in Equation 1.1, restricting the sample to those periods of the game when reports were *not* provided, and including a dummy for whether the participant produced at least one report in the YES category ( $Understand_i$ ). The results showed that participants who understood FI were 22% more likely to exhibit FI-like behavior in rounds without reports ( $P < 0.001$ ; Table 5, column 2). When we controlled for the content of reports, including a variable for the number of reports provided in each category, the coefficient on  $Understand_i$  lost significance, but the coefficient on the number of YES reports was significant and positive, suggesting that each additional report in the YES category made FI-like behavior 11% more likely ( $P < 0.05$ ), while a NO report decreased its likelihood by 12% ( $P < 0.01$ ; Table 5, column 3). We take these results to mean that how well a participant understands FI is a better predictor of behavior than whether or not he or she understands it at all.

Our final comments in this section deal with learning. Recall that while FI-like behavior became more likely with time (Table 2), the occurrence of YES reports did not (Table 3, panels b and c). To reconcile these results, we interacted  $t$  and  $vr_{it}$  with the  $Understand_i$  dummy in the specification described by Equation 1.1,

---

<sup>15</sup>Since some subjects failed to provide reports, the report dummies are not collinear. The no report case serves as a baseline.



	(1)	(2)	(3)	(4)
	Periods with reports	Periods without reports	All periods	All periods
Dependent variable=FI behavior				
Period no.	0.006 (0.007)	0.007**** (0.002)	0.007**** (0.002)	
2.Treatment	-0.066 (0.096)	-0.081 (0.063)	-0.102* (0.059)	-0.080 (0.064)
3.Treatment	0.023 (0.092)	-0.078 (0.080)	-0.094 (0.095)	-0.098 (0.069)
5.Treatment	-0.110 (0.079)	-0.070 (0.066)	-0.066 (0.090)	-0.097 (0.062)
YES report	0.359**** (0.076)			
MAYBE report	0.190** (0.079)			
NO report	-0.124 (0.086)			
Understanding of FI		0.221**** (0.048)	-0.093 (0.092)	
No. of YES reports			0.112** (0.053)	
No. of MAYBE reports			0.019 (0.041)	
No. of NO reports			-0.120*** (0.045)	
Period no. if FI understood				0.017**** (0.003)
Period no. if FI not understood				-0.001 (0.003)
Effect of answering ques. if FI understood				0.144**** (0.032)
Effect of answering ques. if FI not understood				0.030 (0.026)
Observations	446	2202	2202	2648

Table 5: **Participants' thinking accounted for in panel analysis of behavior.** Marginal effects of logit regressions. Subject-clustered standard errors in parentheses. \*  $P < 0.1$ , \*\*  $P < 0.05$ , \*\*\*  $P < 0.01$ , \*\*\*\*  $P < 0.001$ . The first two columns only use periods of the game without reports. Understanding of the FI argument makes FI-like behavior more likely (columns 1-3). Participants who understand FI are more likely to exhibit learning than their peers (column 4).

allowing learning and the response to  $vr_{it}$  to depend on whether FI was understood by the participant. The results ( Table 5, column 4) showed that only participants with an understanding of FI exhibited significant learning ( $P < 0.001$  for subjects who understood FI vs.  $P = 0.7$  for those who didn't) and changed their behaviors in rounds with questions ( $P < 0.001$  vs.  $P = 0.241$ ). We interpret our findings to mean that while FI thinking is unlikely to be learned in the course of an experimental session, participants who have a preliminary understanding of the FI logic can be encouraged with time and by being asked to consider the game more carefully.

## 1.5 Discussion and Robustness

Because the game was played repeatedly with multiple partners, it is possible that participants learned FI thinking from players they were matched with. As mentioned above, however, feedback had no measurable effect on participants' thinking in three out of four of our treatments (i.e., for 170 out of 230 participants). We take this fact to indicate that subjects were unlikely to have learned FI from others, justifying the use of subject-clustered standard errors in our regression analysis. Our results, however, were robust to other model specifications. Thus, using random or fixed effects models, clustering the standard errors by session, or including session fixed effects produced the same basic findings.

Some comments on our belief elicitation procedure are in order. First, when one player is making a guess about another's thoughts, only the former individual gets paid in the case of a match. The other player's report, in this case, is either not incentivized (if own behavior is explained), or incentivized by paying this player more if his guess of the partner's behavior is closer to the truth. Our design therefore does not induce a coordination game with multiple equilibria in statements, where a player has an incentive to state a less sophisticated explanation than his or her true belief if he thinks his matched player might think in a simpler way. The player would only want to do this if he or she is altruistic. Secondly, as noted by Blanco et al. (2010), an experiment that elicits participants' beliefs about their partners' behaviors is subject to a potential hedging problem. In our design, a subject was paid for guessing his partner's behavior only in the case the subgame is played. In this situation, we find that first movers play a risk-neutral best response to their beliefs 82% and second movers 83% of the time. We conclude that most of the reports show no evidence of hedging.

Evaluators were asked to judge the quality of a strategic thinking, so the quality of their performance is essential. To test it, we compared it with our own evaluations of reports in the longer treatments. We found a high degree of correlation between our ratings of first movers' reports and those provided by the evaluators ( $r = 0.86$ ,  $P < 0.001$ ). While this high correlation is not surprising, given that we instructed the evaluators how to rate the reports, it serves as a good test of the evaluators' performance. Note that there is little reason to be concerned about experimenter demand effects. First, we take our definitions and examples of what constitutes forward induction thinking to be uncontroversial. Second, the evaluators were given no reason to believe that we, the experimenters, were interested in finding that many participants exhibit forward induction reasoning, or that few participants do. Indeed, what we told the evaluators is that we are interested in their true beliefs (see footnote five).

As mentioned in [Section 1.3](#), understanding of FI requires a higher degree of strategic sophistication by the first mover than by the second. Our own categorizations, which considered reports of players in both roles, allowed us to address how reports of first and second movers relate. We find that 61 subjects in T1 and T2 failed to provide a YES report. Of the remaining subjects, 22 provided YES reports in both roles; 12 provided it as second but not first movers, and 13 provided it only in the role of first movers. We also observe that *in a given period* a first mover's report is more likely to display evidence of FI than that of a player in the opposite role (23% vs. 37% of reports; distributions of ratings are different with  $P < 0.05$ ). Thus, while the intuition that it is easier to display evidence of FI as a second mover is confirmed, if subjects are given multiple opportunities to report on their thinking, they are equally able to show evidence of FI in either role.

## 1.6 Conclusion

We find strong evidence that FI logic guides thinking and behavior of a substantial number of experimental participants. The display of this ability is in part influenced by the activity of analyzing and explaining one's thought process, and in part by the observation, although limited, of the actions of others: the more often a subject observes their partners choose *In*, the more likely he is to exhibit FI thinking. There is also a link between understanding FI and behaving according to its logic: subjects that show evidence of SBR in their reports are more likely to exhibit FI

behavior. Thus, FI, and not other factors, such as first mover advantage, induces largest frequencies of behavior predicted by the theory.

Previous experimental research had provided only mixed evidence in support of the FI hypothesis. Our results, based on a different method which provides richer information suggest instead that FI thinking models accurately the reasoning process of a significant amount of participants. This confirms existing findings on the specifically human ability to make predictions by pure reasoning (Teglas et al., 2011). In the game we used, the second mover observes information on their opponent's first stage behavior, and, guided by an abstract assumption on players' behavior (SBR), makes a prediction about a novel situation (the second stage of the game). It is likely that previous experimental studies of FI did not provide supportive evidence because FI is not easily identified by behavior in a single game: using only behavioral data, it is impossible to quantify how much the effect of the focal point contributes to ((*In*, *Top*), *Left*) being played. Our methods circumvent this problem by treating FI as a cognitive process and looking for it in participants' explanations of their thoughts and behavior, rather than their selected actions.

Finally, a fraction of one third of players is substantial but is a minority. Our subjects, however, were inexperienced, and were incentivized with amounts typical of laboratory economics experiments. Experienced or professional players, in real life environments, are likely to understand and use the forward induction logic more frequently. This hypothesis should be tested in future research.

## 2 Personality and Bargaining Power<sup>16</sup>

### 2.1 Introduction

In a perfectly competitive economy, imputations of value to factors of production are completely determined (Clark, 1902). Without perfect competition, however, market forces alone fail to determine these imputations (Edgeworth, 1881). In the context of organizations, Knight (1921, II.IV.4) reconciles this indeterminacy as follows:

“There are many productive organizations consisting of small numbers of rather unique agents which very effectively supplement each other and are not so effectively demanded elsewhere. In such a case competition does not afford means of distributing the entire yield of the group among its members; an appreciable part of it resists automatic division and remains a joint product, dependent on the peculiar effectiveness of the particular organization. Many partnerships illustrate this point. Imputation goes as far as the group, giving that its proper income, but fails to distribute accurately within it. In case of a partnership this division between the members is usually made on ethical grounds or on the basis of ‘bargaining power,’ sheer personal force. In industry at large the special product of the organization above that competitively assigned to its components is likely to go, largely at least, to the entrepreneur, though bargaining power or the strategic situation always plays a large part in the proceedings.”

Knight’s argument involves two steps. First, in many organizations, market forces fail to determine its members’ individual imputations. Secondly, this indeterminacy is often resolved by *personal force*. The first step epitomizes an extensive literature on team production, from Alchian and Demsetz (1972) and Holmstrom (1982) to Prendergast (1999), Levin (2003) and beyond. In this paper, we take as given this first step and focus on the second step of Knight’s argument. We offer a formal, quantifiable interpretation of personal force in terms of psychological factors—specifically, individual personality traits—and experimentally test Knight’s hypothesis. Our results show that personal force plays a significant role in the imputation of value.

Individual personality traits can be thought of as enduring behavioral patterns and responses to environmental cues (Almlund et al., 2011). The classic “Big Five”

---

<sup>16</sup>Joint work with David Rahman.

framework (Costa and McCrae, 1992), based on respondents' answers to questionnaires, measures personality along several dimensions: openness to new experience, conscientiousness, extraversion, agreeableness and neuroticism. This measurement framework is widely accepted by psychologists for many reasons, amongst them robustness (Goldberg, 1993), its strong relationship with relevant configurations in the brain (DeYoung et al., 2010), and its ability to predict individual outcomes and choices. Thus, personality is related to long-term individual characteristics such as income, education, health and relationship status (Borghans et al., 2008, and references therein). Recently, interdependencies between personality traits and economic preferences have also been documented (Anderson et al., 2011; Becker et al., 2012). Overall, there is general consensus in the literature that individuals' economic choices and outcomes are significantly affected by their own personality.

We, on the other hand, designed an experiment to test a different hypothesis—that individuals' economic choices and outcomes are significantly influenced by personality traits of *others*. The experiment proceeded roughly as follows.<sup>17</sup> First, each subject completed a Big Five personality questionnaire. The participants were then randomly matched into hierarchical two-person teams, consisting of a worker and a manager, whose members interacted over several periods. Workers performed the same repetitive task every period, which we used to measure and control for productivity, and which translated stochastically into monetary earnings that accrued to the manager. The worker's remuneration was solely the manager's decision, and it came at the manager's own expense. At the end of each interaction, team members completed a personality questionnaire on behalf of their partners, which gave us a measure of perceived personality traits. Subjects were then randomly re-matched.

Our experimental design differs from most in the bargaining literature in three respects. One is that we allowed subjects to regularly interact face-to-face and engage in free-form communication, giving them the opportunity to gradually absorb each others' personality traits. Second, we introduced a reasonable amount of uncertainty in the experiment. This obscured behavioral prescriptions based on ethical grounds and opened the door for greater variation, including personality-driven variation, in behavioral patterns. Third, subjects' reported perceptions of each other's traits provided instruments for identifying a causal effect of personality.

---

<sup>17</sup>See Section 3.3 for further discussion and justification of our design decisions.

Although managers on average allocated 50% of output to workers, our experiment was successful in producing substantial variation in inequality (Table 6). Moreover, this variation could not be explained by variation in output (Table 7). There was, however, significant correlation between the worker’s experienced inequality and her personality traits, which persisted even when physical characteristics were controlled for (Table 8). To assess the causal effect of the worker’s personality, we focused on agreeableness, using the worker’s evaluation of her other manager as an instrument. We found the effect of agreeableness to be statistically and economically significant (Table 20): An increase in a worker’s agreeableness by one standard deviation led the manager to increase the worker’s income by about 16%. Recognizing that teams interacted over several rounds, we also studied the dynamic effects of personality. Early in the interaction, agreeableness had no significant effect on earnings, but its effect increased progressively over time to achieve overall significance (Figure 4). This result is consistent with the hypothesis that subjects gradually learned each others’ personality traits as the experiment proceeded. A priori, it seems plausible that managers might pay agreeable workers less because they would tend to be more accepting of harsher terms. On the other hand, managers may also be inclined to reward workers with a higher opinion of them (as in the psychological games of Geanakoplos et al, 1989). We found that agreeable workers had significantly less favorable opinions of their managers, which suggests that the main channel through which agreeableness translates into lower earnings is the former one: Managers found agreeable workers more docile and decided to pay them less (Table 10).

Our results are important for understanding the psychological sources of bargaining power and, more generally, influence. First, bargaining is a basic facet of economic activity, yet the sources of comparative bargaining advantage do not seem to be well articulated in economic theory. Cooperative solutions, such as Nash bargaining (Nash Jr, 1950) and related variants, take it as given, and noncooperative solutions have so far been unable to usefully incorporate psychological factors. Rubinstein (1982) offers impatience and institutional details (temporal monopoly) to explain bargaining power,<sup>18</sup> yet neither of these issues is practically relevant in our experiment or many real-world situations. Similarly, the model of Abreu and Gul (2000), based on reputation, provides a language with which to express differences in bargaining outcomes, but no guidance whatsoever for the determinants of such reputation.

---

<sup>18</sup>A related extension due to Binmore et al. (1986) adds risk aversion as a possible explanation, provided certain institutional assumptions are met (e.g., a random deadline).

Secondly, even if there was an accepted theoretical relationship between psychological traits and bargaining power, it could only deliver qualitative predictions. In order to *measure* this effect quantitatively, it is necessary to explore the issue empirically. The econometric studies of Seibert and Kraimer (2001), Heckman et al. (2006) and others have the drawback that they study long-run incomes without being able to distinguish between influence or bargaining power and productivity in any specific situation, let alone disentangle the relative values of interpersonal traits and performance. This problem motivates an experimental approach to improve our understanding of just how people’s psychology contributes to their income.

Although there is a vast experimental literature on bargaining, as well as some relating bargaining and personality, it is unable to address our main hypothesis. Most of this literature attempts to explain an individuals’ propensity to share as a function only of their own personality (Brandstätter and Königstein, 2001; Ben-Ner and Kramer, 2011; Anderson et al, 2011). In fact, in these experiments there was no possibility for subjects to learn the personality of counter parties, as interactions were either hypothetical or anonymous. We, however, are interested in measuring the effect of one party’s personality on another’s decision. To accommodate this possibility, our design allowed subjects to learn each others’ personality traits by giving them the opportunity to regularly interact face-to-face and communicate freely.<sup>19</sup>

An exception to this literature is the work of Morris et al. (1999), who analyzed an experiment where MBA students bargained face-to-face over mock salaries. There are important differences between their work and ours in terms of both method and focus—we discuss them at length in Section 2.2 below. In summary, they framed their experiment in a way that reduced the relevance of actual personality, and they focused on understanding how bargaining outcomes and behavior biased perception of personality, rather than the effect of personality on bargaining outcomes.

## 2.2 Literature on Income and Personality

This paper is motivated partly by a well-documented relationship between income and personality. Heckman et al. (2006) estimate a wage equation that significantly

---

<sup>19</sup>Free-form face-to-face interaction is standard practice in psychology and organizational behavior (Thompson et al, 2010, and references therein). Face-to-face communication is less common in economics, but accepted (e.g., Mobius and Rosenblat, 2006). See Section 3.3 for further discussion.



relates earnings with cognitive skills (such as IQ) and noncognitive skills (such as personality), suggesting that (p. 1) “[...] personality traits, persistence, motivation and charm matter for success in life.” A number of studies looking at effects of individual personality traits identified a negative relationship between agreeableness and income for both men and women (Mueller and Plug, 2006; Nyhus and Pons, 2005; Ng et al, 2005; Rode et al, 2008). However, these analyses leave out important details about how a worker’s personality affects his or her income, as well as the role of others in wage determination. In other words, they fall short of being able to explain just how “charm” (for instance) matters for success in life. Thus, it cannot be inferred from Heckman et al.’s wage equation whether personality increases wages because it motivates individuals towards more productive behavior or more rent-seeking behavior, such as bargaining skills, which may be unproductive, as Knight (1921) suggests. One goal of our study is to disentangle quantitatively these different motivations in a richer model of wage determination, thus beginning to open the “black box” behind the relationship between personality and earnings.

Personality also has a well-documented effect on economic preferences (e.g., Borghans et al, 2008; Becker et al, 2012). More relevant to our paper’s main results is perhaps the observation that agreeable people are more altruistic in dictator (Ben-Ner et al, 2008) and trust (Anderson et al, 2011) games. We should emphasize, though, that our experiment differs from others in the personality and bargaining literature by virtue of focusing on the link between one’s decisions and personality traits of *other people*. Hence, our paper is closer in spirit to the work of Ben-Ner and Kramer (2011), who studied the effect of kinship on the amount received in a dictator game, and Judge et al. (2012), who found that one’s personality—particularly agreeableness—affected another’s estimate of job growth potential. Both of these studies, however, used *hypothetical* descriptions of people as explanatory variables. For our purposes, real, direct interaction was important to allow personality traits to both express themselves endogenously and translate into bargaining power, rather than be communicated exogenously. Incentives were hypothetical in these studies, too. This is important, as according to Camerer and Hogarth (1999), excluding financial incentives may increase certain behavioral traits associated with personality, such as generosity and risk-seeking. This motivates our use of monetary transactions to clarify the relation between surplus division and personality.

Arguably, the study closest to ours in method is Morris et al. (1999), which also used face-to-face interaction in a bargaining environment. This study, however, fo-

cused on how bargaining outcomes and behavior biased perceptions of personality. In particular, the authors did not measure how personality traits of other people affect bargaining decisions. They also argued that the behavior of participants in their experiment was mostly driven by “situational” rather than personality factors.<sup>20</sup> This is, perhaps, not surprising, in light of [Morris et al.’s](#) experimental design. As the authors state in their paper (p. 56), “[p]articipants were familiar from negotiation class with the concepts of the value and risk of an alternative option and had been taught guidelines for estimating these from an opponent’s negotiation behavior.” As [Thompson \(1990\)](#) and [Monson et al. \(1982\)](#) argue, personality is more likely to matter when strong behavioral prescriptions, such as those taught to the MBA students in [Morris et al.’s](#) study, are absent. We designed our experiment with this in mind.

### 2.3 Experimental Design

Our motivation for the experiment was to create an environment that resembled the spirit of Knight’s argument and allowed us to test his hypothesis. We matched subjects into teams of two, motivated by the observation that individuals often interact in small groups ([Burke, 2003](#)). By design, the teams did not interact with one another, so there was no competition for team members. This feature of the experiment kept it aligned with [Knight’s \(1921, II.IV.4\)](#) observation that “[t]here are many productive organizations consisting of small numbers of rather unique agents which very effectively supplement each other and are not so effectively demanded elsewhere.” As a result, the division of surplus amongst team members became indeterminate and open to bargaining, and, hence, possibly personal force.

We framed the experiment around a hierarchical organization whose members performed different tasks,<sup>21</sup> to avoid a situation that might easily lead subjects to agree on equal surplus division. This issue is well-documented in experiments, espe-

---

<sup>20</sup>Specifically (p. 53), “[...] important components of bargaining behavior [...] are greatly determined by the economic incentives and constraints a player faces and little determined by personality traits ([Thompson, 1990](#)).” However, [Thompson \(1990\)](#) is much more cautious, admitting that (p. 520) “[...] this conclusion is incomplete and overly simplistic.” Amongst several reasons for this view, she reports that (pp. 520-521) “[Monson et al. \(1982\)](#) suggested that personality is more predictive of behavior in ambiguous situations than in settings in which there are strong prescriptions for behavior.”

<sup>21</sup>Notice, however, that—as seen from the experiment’s instructions ([Appendix 4.9](#))—no explicit hierarchical descriptions of player roles, such as “worker” or “manager,” were imposed on subjects.

cially ones without anonymity. Thus, [Bohnet and Frey \(1999\)](#) show that removing anonymity in dictator games eliminates most variation in offers around equal division. On the other hand, we viewed face-to-face interaction as an important aspect of our experimental design, since it allowed subjects to learn each other’s personality traits. By itself, this element of our design might substantially reduce the variation in offers. For our purposes, however, variation was important to be able to trace the relationship between personality and offers, since without variation there could be none due to personality. Therefore, to compensate for the loss of variation in offers due to lack of anonymity, we subjected team members’ interaction to a reasonable amount of ambiguity and complexity, on the grounds that more ambiguity would give subjects moral “wobble room” for their decisions. This intuition was substantiated experimentally by [Dana et al. \(2007\)](#), who showed that (see their abstract) “[...] fairness decreases substantially when the connection between choices and outcomes is obfuscated.”

Some economists have expressed concern regarding face-to-face interaction in experiments. One reason may be that ([Crawford, 1998](#), p. 293) “[n]onpecuniary influences on preferences are usually suppressed by avoiding face-to-face or nonanonymous interactions [...]” However, these are precisely the influences we are trying to capture. A particularly appealing reason for choosing face-to-face interaction rather than chat messages, phone-based or other types of controlled communication is perhaps best articulated by [Nadler and Shestowsky \(2006](#), p. 165): “[...] when the structure of the negotiation is a complex, potentially integrative negotiation that requires reciprocal information sharing, the inability to see or hear the other person in conjunction with lack of co-temporality can exacerbate initial distrust, leading to reluctance to engage in the kind of reciprocal exchange of information required to reach a high-quality agreement, or any agreement at all, for that matter.” Our environment, described below, is complex enough that this was a potential concern.

Each experimental session was divided into two halves. Subjects were randomly rematched from one half to the next, with subject roles unchanged, so workers remained workers. Perceived personality traits were recorded at the end of each half. Designing our experiment to have these two halves was particularly useful for two reasons. First, it gave us some variation in outcomes for each subject, improving the statistical properties of our sample. Second, it delivered a useful instrument to identify a causal relationship between earnings and endogenous variables. In principle, a worker’s personality may be correlated with other factors unobservable to us that

contributed to the manager’s determination of the worker’s income. We found that a worker’s personality was correlated with her perception of her manager. Since each worker was matched with two different managers, to identify the effect of a worker’s personality on a manager’s remuneration decision, we used the worker’s perception of the other manager’s personality as an instrument. See [Section 2.4.3](#) for details.

### 2.3.1 Details of the Experimental Design

The experiment was programmed and conducted with the software z-Tree ([Fischbacher, 2007a](#)) in the Anderson Hall Social and Behavioral Sciences Laboratory at the University of Minnesota. After completing the Big Five personality questionnaire of [DeYoung et al. \(2007\)](#), subjects familiarized themselves with instructions provided.<sup>22</sup> They were then randomly matched into teams of two, and each team member was randomly allocated the role of worker or manager. Worker and manager sat next to one another in separate carrels and interacted for 15 rounds. Everyone was told that they were sitting next to their teammate after being matched.<sup>23</sup>

In each round, the worker’s job was to complete a repetitive task, borrowed from [Gill and Prowse \(2012\)](#): to move as many sliders as possible, from a total of 24, within an allotted time of 40 seconds. A monetary prize of \$4 was contained behind one and only one of the sliders. Moving a slider meant physically dragging it to position 50 (out of 100) with a mouse. For every slider not moved to position 50, a penny was added to worker’s “penny” account, which was kept separate from the account the manager used to pay the worker.<sup>24</sup> There was therefore a real as well as a monetary cost of effort. We hoped that emphasizing the monetary cost would make it clearer to the managers that workers need to be incentivized. The worker was never informed of whether or not she discovered a prize.

The manager started out with \$5, and had to pay 40 cents in every period, in order to continue the experiment. If and only if the worker discovered a prize, \$4 were added to the manager’s personal account. There was therefore a possibility of the

---

<sup>22</sup>See [Appendix 4.9](#) for the instructions and [Appendix 4.7](#) for the questionnaire.

<sup>23</sup>It is therefore possible that personality had an effect through first impressions even before the subjects were told to talk to each other. E.g., it was shown by [Willis and Todorov \(2006\)](#) that people are able to form first impressions within 100 milliseconds of exposure to a face. The analysis of [Section 2.4.3](#) explores the possibility that the effect of personality changed over time.

<sup>24</sup>This penny was added even if a slider was moved to position 49.

team going bankrupt after 12 periods, in case that no prizes at all were discovered. After observing how many “Top” sliders (i.e., sliders 1 through 12) and “Bottom” sliders (i.e., sliders 13 through 24) the worker moved to position 50, as well as whether or not the prize was found, the manager decided how much to pay the worker.<sup>25</sup> This payment could be any number of cents up to the amount of money the manager accumulated so far. Thus, all of the manager’s start-up funds could be allocated to the worker in the first period, terminating the experiment (because no money is left to continue). On the other extreme, the manager could refrain from paying the worker anything until the very last period. Crucially, decisions of the manager were reversible: any money allocated to the worker by the manager (hence, excluding the worker’s earnings from unadjusted sliders) could be taken back in a subsequent period. Thus, the interaction mirrored a dictator game in that the manager could appropriate the total surplus (minus one dollar, since the manager started out with \$5 and had to pay 40 cents in every period, including the first one) in the very last period. After paying the worker, the manager decided what subset of sliders (Top or Bottom) to recommend to the worker.

The location of the prize-winning slider changed pseudo-randomly according to a Markov process with 75% transition probability for the state (whether the prize was behind a Top slider or Bottom slider) being the same, although the subjects were not informed of this.<sup>26</sup> Conditional on the prize-winning slider a Top slider or a Bottom slider, its location amongst the Top or Bottom sliders was otherwise determined with equal probability of 1/12. Whether the prize was behind a Top slider or Bottom slider was a common event for every team, but the location of the prize within the Top or Bottom sliders was identically and independently distributed across teams.

Every five rounds, the teammates were allowed to talk, face-to-face, for three minutes. Their instructions encouraged discussing the experimental task, but interactions were otherwise unstructured. At the end of the match, subjects were asked to complete a personality questionnaire on behalf of their partner. This concluded the first half of the experiment. For the second half of the experiment, subjects were randomly re-matched with player roles unchanged, so workers remained workers, and

---

<sup>25</sup>The manager had unlimited time to make all of her decisions.

<sup>26</sup>The instructions provided subjects with the following information (“Person A” corresponds to the worker and “Person B” to the manager): “Whether the prize is behind a TOP/BOTTOM slider in the next round only depends on where the prize was in this round. Person A will never know where the prize is. At the end of every round, Person B will see whether or not the prize was discovered. He/she will use this information to make recommendations to Person A.”

the interaction described above was repeated.<sup>27</sup>

## 2.4 Results

172 subjects participated in eight experimental sessions, with session sizes ranging from 10 to 26 subjects. Because each subject took part in two teams, each having two members, this produced data for 172 matches. The remainder of this section is devoted to our main result. A detailed description of the dataset as well as its summary statistics can be found in [Appendix 4.4](#).

### 2.4.1 Income, Inequality and Incentives

We focus on end-of-match outcomes in our statistical analysis. The variables of interest are described below.

- $Effort_{it}$ : the number of sliders adjusted correctly by worker  $i$  in match  $t \in \{1, 2\}$ .
- $Income_{it}$ : worker  $i$ 's total earnings (including pennies for unadjusted sliders) in match  $t$ , measured in dollars.
- $Inequality_{it}$ : the difference between the earnings of worker  $i$  and the earnings of her manager in match  $t$ , measured in dollars.
- $Output_{it}$ : the number of prizes discovered by worker  $i$ 's team in match  $t$ .
- $Recommendations_{it}$ : the number of good recommendations made by  $i$ 's manager in match  $t$ .

The summary statistics of these variable are provided in [Table 6](#). The median difference in earnings is nearly zero, suggesting that managers did keep fairness in mind when deciding how much the worker should be rewarded. The experiment, however, was successful in producing substantial variation in inequality (SD=5.975).

---

<sup>27</sup>The locations of the prize-winning slider were {Top, Top, Top, Bottom, Top, Bottom, Bottom, Bottom, Bottom, Bottom, Bottom, Top, Bottom, Top} in the first half of the experiment and {Bottom, Bottom, Bottom, Top, Top, Bottom, Top, Top, Top, Top, Bottom, Bottom, Bottom, Bottom, Bottom} in the second half. The two halves were otherwise identical in design.

Although subjects' interactions spanned two halves of the experiment, with a different teammate in each, the distributions of the income, inequality, and output variables did not differ across halves ( $P = 0.433$  for income,  $P = 0.141$  for inequality, and  $P = 0.923$  for output according to a Kolmogorov-Smirnov test). Workers moved more sliders in their second match, but the managers made worse recommendations, and similar numbers of prizes were discovered. We discuss this in more detail in the data appendix.

	Mean	SD	Min.	Max.	Median	P-value of K-S test
Effort	99.42	25.03	19	170	100	0.001
Income	8.396	3.952	2.330	19.91	8.100	0.433
Inequality	-1.608	5.975	-18.91	14.55	-0.0550	0.141
Output	4.199	1.732	1	9	4	0.923
Recommendations	7.038	1.880	2	12	7	0.025
$N$	156					

Table 6: Summary statistics of the variables used in the analysis. For each variable, the rightmost column reports the P-value of a Kolmogorov-Smirnov test assessing equality of distributions in the first and second half of the experiment.

Recall that a team could go bankrupt if no prize was discovered for 12 periods, or if the manager did not leave herself enough money to continue to the next period because too much had been allocated to the worker (e.g., the manager may not have understood the instructions). Nine out of the 172 matches were confronted with the former situation, and seven failed to find a prize and become bankrupt as a result. All bankrupt matches were excluded from our subsequent analysis.

To identify how managers allocated output between themselves and their workers, we estimated the following models:

$$Income_{it} = \alpha + \beta Output_{it} + \epsilon_{it} \quad (2.1)$$

$$Inequality_{it} = \alpha + \beta Output_{it} + \epsilon_{it} \quad (2.2)$$

These models suffer from a possible endogeneity issue because we can't rule out that the worker's effort is rewarded (i.e., included in  $\epsilon$ ), and workers are likely to respond to higher incomes by working more. As an instrument for output, we used the number of good recommendations given by the manager. This instrument is valid:

	Income	Inequality
Output	1.412*** (0.470)	-1.267 (0.947)
Constant	2.466 (1.864)	3.714 (3.763)
F-statistic (first stage)	21.61	21.61
F-statistic (second stage)	8.071	1.602
Underidentification test	156	156

Session-clustered standard errors in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ , \*\*\*\*  $p < 0.001$

Table 7: Managers rewarded workers for producing output. Note that while output had a significant effect on workers' incomes, it had no significant effect on inequality.

firstly, because better recommendations lead to more prizes being discovered, so the instrument is correlated with output, and secondly because there is no reason to think that how much the manager decides to pay the worker, conditional on output, is correlated with how good the manager's recommendations are. The reason for this is that the number of good recommendations made by the manager is censored: The manager only observes output, and never finds out how many good recommendations were made when a prize was not discovered (unless, of course, a manager found the prize in every period of a match, but this never occurred).

We estimated the model with two stage least squares, clustering the standard errors by session.<sup>28</sup> The results, reported in Table 7, show that the constant term in this regression did not differ significantly from zero ( $P = 0.324$ ), confirming the observation that managers split the output evenly on average. While output had a significant effect on the worker's income ( $P < 0.01$ ), it had no significant effect on the inequality variable. I.e., variation in inequality could not be explained by variation in output. We show below, however, that a significant portion of this variation could be explained by the worker's personality.

<sup>28</sup>We used session-robust standard errors in all of our statistical analysis.



## 2.4.2 Correlations With Personality

We first assessed correlations of personality with income and inequality, using models of the following form:

$$Income_{it} = \alpha + \beta\Psi_{it} + \gamma\Phi_{it} + \epsilon_i \quad (2.3)$$

$$Inequality_{it} = \alpha + \beta\Psi_{it} + \gamma\Phi_{it} + \epsilon_i \quad (2.4)$$

$\Psi_{it}$  is a vector of  $i$ 's personality characteristics, as well as the personality characteristics of her manager in match  $t$ , and  $\Phi_i$  a vector of additional covariates, such as output and physical traits (beauty and gender). Personality characteristics were  $z$ -scored in our analysis; i.e., each trait had the sample average subtracted, and the difference was divided by the standard deviation. A description of the beauty and gender data is included in [Appendix 4.4](#). Eight sessions are used in the regressions where physical traits are included as controls. As discussed in the appendix, time constraints prevented us from obtaining these data in the first two sessions, although the sessions were identical to the others in all other respects. We considered several specifications of the models described above in addition to the baseline regressions with  $\gamma = 0$ : (1) controlling for output, (2) controlling for physical traits of the worker and manager, (3) controlling for output as well as physical traits. Note that  $Inequality_{it} = Output_{it} - Income_{it}$ . Therefore, using inequality as a dependent variable *without* accounting for output is equivalent to using income as a dependent variable and restricting the coefficient on output to be 1/2. This captures the hypothesis that the manager “splits the pie” equally on average, while deviations from this norm reflect influences of other factors such as the worker’s personality. When output is included as a control, we report the results of the income regression.

[Table 8](#) shows the results of two joint hypothesis tests for each of the considered models. The first hypothesis is that the effect of every personality trait of the worker is zero. The second is an analogous hypothesis concerning the effect of the manager’s personality. The results, collected in [Table 8](#), show that a worker’s personality traits had a jointly significant effect on her earnings in every specification (largest  $P = 0.0201$ ), while the joint effect of a manager’s traits on her income or experienced inequality was always insignificant (smallest  $P = 0.1222$ ). The finding that the manager’s personality had no significant effect on income or inequality is consistent with [Morris et al. \(1999\)](#), who argue that bargaining outcomes are often determined by “situational” factors rather than own personality characteristics. [Morris et al. \(1999\)](#), however, did not measure the effect of other people’s personality traits on

Additional controls	Personality of worker	Personality of manager	Sessions	$N$
None	$p = 0.0095$	$p = 0.1415$	10	156
Output	$p = 0.0105$	$p = 0.1895$	10	156
Physical traits	$p = 0.0022$	$p = 0.6333$	8	123
Output and physical traits	$p = 0.0044$	$p = 0.4834$	8	123

Additional controls	Personality of worker	Personality of manager	Sessions	$N$
None	$p = 0.0134$	$p = 0.5115$	10	156
Physical traits	$p = 0.0201$	$p = 0.1222$	8	123

Table 8: Correlations between personality and the worker’s income/average experienced inequality.  $P$  statistics test the null hypothesis that the five personality traits of the worker and manager, respectively, are jointly significant, after controlling for (i) nothing else, (ii) output, (iii) physical traits (gender and beauty), and (iv) both output and physical traits.

own bargaining decisions. To our knowledge, the finding that this effect is significant has not been reported in the literature.

### 2.4.3 Identification

Correlations between personality traits and bargaining outcomes leave unsettled the issue of causality. We used the variation in workers’ ratings of their managers, which were related to workers’ agreeableness, to identify a causal effect of agreeableness on earnings. Specifically, we estimated the following models of manager  $t$ ’s payment decision, using worker  $i$ ’s evaluation of her *other* manager  $t'$  as an instrument for  $Agreeableness_i$ :

$$Income_{it} = \alpha + \beta Agreeableness_i + \epsilon_{it} \quad (2.5)$$

$$Inequality_{it} = \alpha + \beta Agreeableness_i + \epsilon_{it} \quad (2.6)$$

The instrument we used requires justification. Recall that workers and managers evaluated each other’s personalities through questionnaires at the end of each interaction. As stated in DeYoung et al. (2007) (p. 883), “All of the positive poles of the Big Five are socially desirable, whereas all of the negative poles are socially undesirable (Neuroticism is reversed [...] and labeled *Emotional Stability*).” We define  $Desirability_{it'}$ ,  $i$ ’s expressed desirability of manager  $t'$ , as the sum of worker  $i$ ’s ratings of this manager’s extraversion, agreeableness, conscientiousness and openness

minus her rating of the manager’s neuroriticism, and use it to instrument for the effect of  $i$ ’s agreeableness on the payment decision of manager  $t$ . This is valid for two reasons. First,  $Desirability_{it'}$  is highly correlated with  $i$ ’s agreeableness ( $\rho = 0.3869$ ,  $P < 0.001$ ). In Section 2.4.4, we study this relationship in more detail and argue that while the worker’s reported perceptions of her manager’s personality are correlated with agreeableness, they are not correlated with other personality traits of the worker. Second, since workers and managers were randomly matched, it is reasonable to assume that  $E[Desirability_{it'} \cdot \epsilon_{it} | Agreeableness_i] = 0$ , where  $t' \neq t$ . In words, this assumption amounts to claiming that—conditional on the worker’s agreeableness—just how a manager pays a given worker is independent of how the worker rates the *other* manager to whom she is matched in the other half of the experiment.

	Income	Inequality	Income	Inequality
Agreeableness	-1.659 (1.081)	-2.523** (0.982)	-1.356** (0.563)	-2.761** (1.187)
Output			1.474*** (0.490)	-1.161 (0.992)
Constant	8.660**** (0.472)	-1.303* (0.680)	2.385 (1.971)	3.640 (4.005)
F-statistic (first stage, agreeableness)	21.10	21.10	20.02	20.02
F-statistic (first stage, output)			18.31	18.31
F-statistic (second stage)	2.106	5.902	6.289	3.150
Underidentification test	0.0205	0.0205	0.0155	0.0155
Observations	154	154	154	154

Standard errors in parentheses  
\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ , \*\*\*\*  $p < 0.001$

Table 9: The effect of the worker’s agreeableness on her bargaining power, with the worker’s rating of the (other) manager used as an instrument.

We estimated the model described above with two stage least squares and session-clustered standard errors. Workers in several teams provided the same answer (“Neither Agree Nor Disagree”) for every question regarding their boss, and these teams were excluded from the sample. The coefficient estimates are reported in Table 20. While the worker’s agreeableness had a significant effect on her experienced inequality ( $P < 0.05$ ), the coefficient in the income regression was not significant ( $P = 0.142$ ). All standard tests of instrument validity were passed by these specifications: the  $P$ -

value of the Kleibergen-Paap underidentification test was equal to 0.0205, and the first stage  $F$ -statistic was 21.10, suggesting that the instrument for agreeableness was not weak.

When we re-estimated the model including output (measured in terms of the number of discovered prizes) as an additional control, and using the manager’s recommendation as an instrument for output, we found the effect of agreeableness to be significant in both specifications ( $P < 0.05$ ). The null hypothesis of underidentification was rejected with  $P = 0.0155$ . The multivariate first stage  $F$ -statistics of Angrist and Pischke (2008) were 20.02 for agreeableness and 18.31 for output.

Notice that the point estimates were quite similar for models with and without output included as a control (-1.659 and -1.356 for income; -2.523 and -2.761 for inequality). These point estimates, considered together with mean earnings of 8.396 reported in Table 6, suggest that the worker’s income decreased by 16%-20% per standard deviation of agreeableness. Thus, the effect of agreeableness on the manager’s decision was economically as well as statistically significant. We highlight this finding below as our first major result:

**Result 1.** THE WORKER’S AGREEABLENESS CAUSED THE WORKER TO BE PAID SIGNIFICANTLY LESS BY HER MANAGER.

If the effect of agreeableness on the worker’s earnings was due to face-to-face interactions with the manager, one may hypothesize that this effect strengthened with the number of interactions as subjects became acquainted with their team members and gradually assimilated their personality traits. To study the dynamics of the effect of agreeableness over time, we estimated the marginal effect of agreeableness on the manager’s decision for each period of the interaction. I.e., we re-estimated the models described above taking income and inequality in periods preceding the very last one. This led to 30 regressions. An additional 30 regressions included output as a control, with the manager’s recommendations serving as an instrument as before.

We plot the coefficients on agreeableness together with 95% confidence intervals in Figure 4. In every panel of the figure, the marginal effect of agreeableness was indistinguishable from zero for the first several periods of the interaction.<sup>29</sup> With time, the

---

<sup>29</sup>One exception is the specification where inequality is used as a dependent variable and output is not controlled for, where the effect of agreeableness was significant at a 10% level in the very first period. In the second period, however, the largest  $P$ -value is equal to 0.631 in any model.

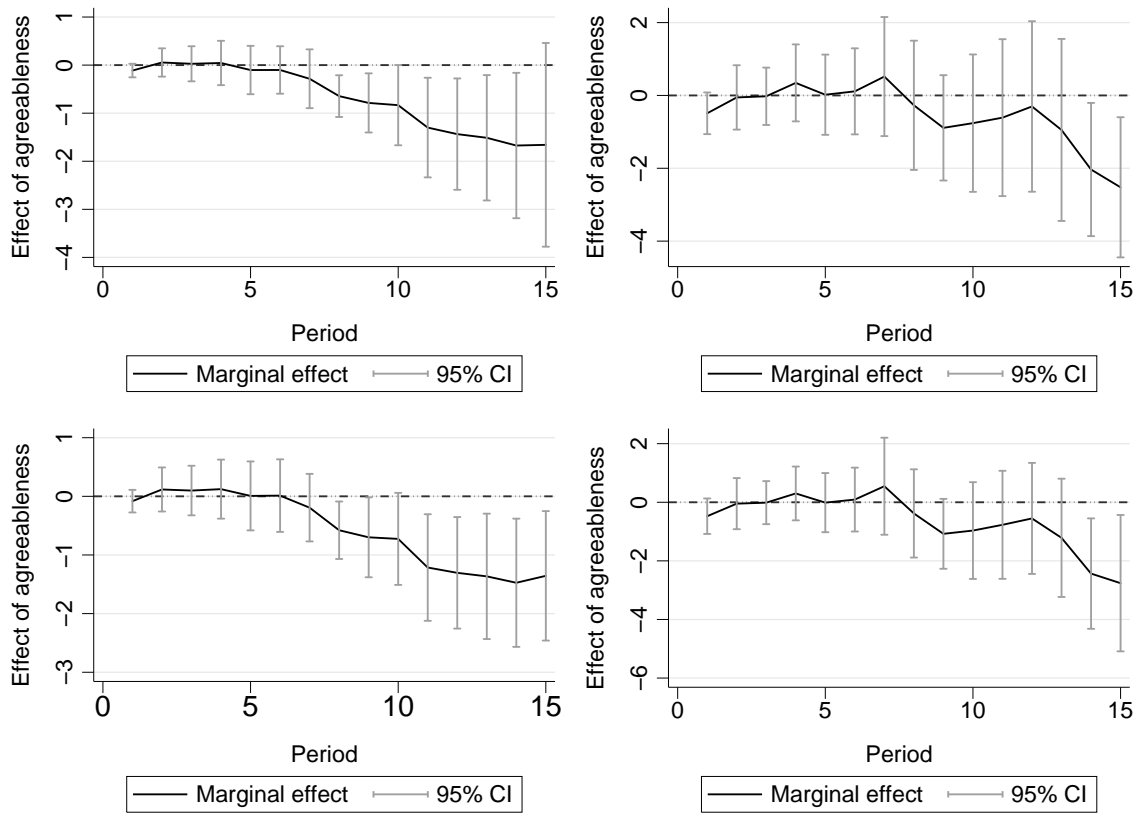


Figure 4: The change in the marginal effect of agreeableness over time. 95% confidence intervals are plotted around the marginal effects at each period of the interaction.

effect grew more negative, and eventually reached significance. When inequality was the dependent variable, it became significant only at the very end of the interaction; when the dependent variable was income, it became significant earlier. Note that even when income serving as the dependent variable and output not controlled for, the effect of agreeableness was significant in several *late* periods of the interaction. We take the finding that the effect of agreeableness became stronger over time to be our second major result, highlighted below.

**Result 2.** THE EFFECT OF AGREEABLENESS WAS INSIGNIFICANT DURING THE BEGINNING OF THE INTERACTION, BUT GREW IN MAGNITUDE AND BECAME SIGNIFICANT OVER TIME.

We interpret this result as evidence for the hypothesis that it took time for the manager to become acquainted with her worker’s personality and respond to it.

#### 2.4.4 Evaluation of Others Through Questionnaires

We now turn to the analysis of subjects’ perceptions of their teammates. Three subjects entered the same answer (“Neither Agree Nor Disagree”) for every item in the survey they filled out on their teammate’s behalf; these subjects were excluded from subsequent analysis. To study the relationship between the worker’s perception of the manager and personality, we estimated regressions of the following sort:

$$Perception_{itk} = \alpha + \beta\Psi_{it} + \epsilon_{it}, \quad (2.7)$$

where  $Perception_{itk}$  stands for  $i$ ’s perception of the  $k^{th}$  trait of manager  $t$ . The results are reported in Table 10. We found that agreeable workers saw their teammates as being more open ( $P < 0.01$ ), more conscientious ( $P < 0.01$ ), more extraverted ( $P < 0.1$ ), less neurotic ( $P < 0.001$ ), and more agreeable ( $P < 0.01$ ). Thus, agreeable workers perceived their managers as being more socially desirable, i.e., liked them more. This is confirmed in the last column of the table, where the worker’s expressed affinity of the manager is regressed against the same explanatory variables. We take this to be our third major finding.

**Result 3.** AGREEABLE WORKERS EXPRESSED MORE POSITIVE RATINGS OF THEIR TEAMMATES’ PERSONALITIES.

Note that we cannot reject the hypothesis that the effect of personality traits other than agreeableness on the worker’s evaluation of the manager was zero. Thus, to the

	Perc. N.	Perc. A.	Perc. C.	Perc. E.	Perc. O.	Desirability
N. of worker	0.0871* (0.0472)	-0.00422 (0.0907)	-0.104*** (0.0276)	-0.0642 (0.0561)	-0.0581 (0.0383)	-0.318 (0.216)
A. of worker	-0.189**** (0.0330)	0.244*** (0.0590)	0.142*** (0.0424)	0.125* (0.0661)	0.144*** (0.0429)	0.844**** (0.154)
C. of worker	0.0294 (0.0433)	0.0227 (0.0937)	-0.0108 (0.0334)	0.00775 (0.0682)	-0.0153 (0.0722)	-0.0251 (0.245)
E. of worker	-0.0325 (0.0508)	0.00323 (0.0775)	0.0637 (0.0479)	0.0415 (0.0369)	0.0970* (0.0492)	0.238 (0.199)
O. of worker	0.0634 (0.0456)	0.0514 (0.0741)	0.0321 (0.0596)	-0.0249 (0.0401)	0.0945** (0.0396)	0.0897 (0.119)
N. of manager	0.0418 (0.0418)	-0.00316 (0.0396)	0.0380 (0.0467)	-0.0488 (0.0429)	0.0389 (0.0410)	-0.0168 (0.104)
A. of manager	0.00343 (0.0403)	0.0397 (0.0587)	-0.0680 (0.0451)	0.0490 (0.0650)	-0.00954 (0.0534)	0.00773 (0.210)
C. of manager	0.0290 (0.0341)	0.0283 (0.0531)	0.0321 (0.0381)	-0.0608 (0.0569)	-0.0252 (0.0297)	-0.0545 (0.118)
E. of manager	-0.0812 (0.0486)	0.00231 (0.0522)	-0.00175 (0.0236)	0.118** (0.0419)	0.0482 (0.0360)	0.248 (0.151)
O. of manager	0.0541 (0.0411)	-0.111* (0.0520)	0.0125 (0.0306)	-0.0865* (0.0439)	0.0219 (0.0348)	-0.217 (0.155)
Constant	2.556**** (0.0334)	3.505**** (0.0479)	3.568**** (0.0527)	3.257**** (0.0745)	3.419**** (0.0624)	11.19**** (0.214)
Observations	154	154	154	154	154	154

Standard errors in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ , \*\*\*\*  $p < 0.001$

Table 10: The worker's perceptions and personality traits of both teammates.

extent that the worker was rewarded for her positive perception of the manager, this was outweighed by a disagreeableness premium. Our interpretation is that managers found agreeable workers more willing to accept harsher terms, and hence paid them less.

#### 2.4.5 Robustness Checks

Our main estimation results are reported in [Table 8](#) and [Table 7](#) under several specifications. These results survived additional robustness checks such as including session and half fixed effects. These robustness checks are collected in [Appendix 4.5](#). The effect of agreeableness on inequality remained significant when we measured inequality as the share of output allocated to the worker and estimated the model with two stage least squares. Such a specification, however, is problematic because this dependent variable takes on values in the unit interval ([Papke and Woolridge, 2008](#)); we therefore take our inequality variable (difference in earnings) to be a more natural unit of analysis. When we used the worker’s reported perceptions of the manager’s five traits as instruments, instead of aggregating them into a desirability variable, the effects of agreeableness on bargaining outcomes remained significant and comparable, although the first stage F-statistics suggested that perceptions of some traits were only weakly correlated with agreeableness once perceptions of the other traits were taken into account. We therefore use our aggregated measure of perception, which has a natural interpretation as an index of social desirability, and which survived all standard tests of instrument validity, as an instrument in the main text.

#### 2.4.6 The Worker’s Effort

The preceding analysis focused on the decisions of the manager. We also studied the worker’s effort decisions, and their relation to personality. We did not find strong evidence to support the hypothesis that either the worker’s or the manager’s personality affected the worker’s effort.

Did workers respond to incentives? To answer this question, we estimated the following model:

$$Effort_{it} = \alpha + \beta Income_{it} + \epsilon_{it}. \tag{2.8}$$



	Effort	Effort	Effort
Income	4.660**** (0.450)		4.479**** (0.595)
Agreeableness		-2.417 (9.739)	4.470 (6.915)
Constant	60.30**** (5.794)	99.74**** (3.339)	61.23**** (5.861)
F-statistic (first stage, income)	87.99		91.67
F-statistic (first stage, agreeableness)		22.55	55.02
F-statistic (second stage)	95.77	0.0551	55.02
Underidentification test	0.0037	0.0192	0.0193
Observations	156	156	156

Standard errors in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ , \*\*\*\*  $p < 0.001$

Table 11: Determinants of the worker’s effort.

Because effort is likely to be rewarded by the manager, this model is subject to a potential endogeneity issue. To address this concern, we used output as an instrument for income. On the one hand, output and income are highly correlated. On the other hand, output was not observed by the worker, and hence should be uncorrelated with the worker’s effort decision conditional on income. We note, however, that the latter assumption is not innocuous. In particular, it is possible that the worker’s effort was influenced by what the manager told the worker, making the latter variable correlated with output. The two stage least squares estimation results, reported in Table 11, demonstrate that higher incomes induced the workers to exert more effort ( $P < 0.001$ ), providing some evidence that the worker’s effort was influenced by her income.

We find no significant correlation between the worker’s effort and her own personality. The results of hypothesis tests analogous to those shown in Table 8 with income and inequality as dependent variables are reported in Table 19.<sup>30</sup> The second

<sup>30</sup>Note, however, that these correlations show that the worker’s effort was correlated with the manager’s personality at a 10% level. When the manager’s agreeableness was included as a covariate in the model described by Equation 2.8—with the manager’s perception of the worker in the other

column of Table 11 reports the results of a regression that adds agreeableness to the model described by Equation 2.8. The results of this regression showed agreeableness to have no significant effect on effort ( $P = 0.518$ ), providing additional evidence that the effect of the worker’s agreeableness on the manager’s decisions cannot be attributed to differences in the worker’s productivity.

Additional controls	Personality of worker	Personality of manager	Sessions	$N$
None	$p = 0.3164$	$p = 0.0745$	10	156
Physical traits	$p = 0.1102$	$p = 0.0684$	8	123

Table 12: Correlations between personality and the worker’s effort.

## 2.5 Conclusion

We studied the effect of personality on bargaining power in a controlled experiment, designed to broadly reflect team production in an organization. The combined results reported in this paper point to three main observations: (i) agreeableness of an individual at the bottom of a given hierarchical relationship is associated with decreased bargaining power, (ii) this effect becomes stronger as the individual’s personality is learned, and (iii) the effect is not due to the fact that agreeable individuals tend to view their superiors more favorably. Agreeableness is typically defined as “the tendency to act in a cooperative, unselfish manner” (Becker et al., 2012, Table A.1). We suggest that this tendency is perceived by managers, and, perhaps subconsciously, exploited.

In the future, it would be interesting to relax the bargaining problem we studied here and understand just how robust our results are to specific details of the economic environment. For instance, although we held the hierarchy fixed in our experiment, some evidence suggests that personality is related with status-seeking behavior (Kyl-Heku and Buss, 1996), and, hence, one’s status too. Therefore, the effect of personality on a typical organization is likely to be much more complex than the one observed in this paper. Nevertheless, an important motivation of this study is to open the door for detailed experimental analysis of personality in environments that are both strategic and not anonymous, as is the case in many important economic relationships. Potential applications of this idea range from a deeper under-

---

half used as an instrument—we found no significant causal effect.

standing of earnings determination in organizations to the relevance of Luxembourg in European politics.

Finally, it seems reasonable to hypothesize that disagreeable workers exhibited characteristic behavioral traits that were effectively unobserved to us as experimenters (speaking in a louder voice, etc.). The goal of this study was to investigate the effect of other people's personality on one's economic decisions, rather than trying to understand in depth the channels through which personality traits express themselves. Understanding these channels more deeply, as well as how they interact with strategic considerations, seems to us an exciting topic for future research.

## 3 Cooperative Institutions<sup>31</sup>

### 3.1 Introduction

Cooperative agreements are often complicated by limitations on how much information is available to the parties involved. Firms in an industry attempting to collude, for example, cannot reliably verify every relevant decision made by their competitors. Likewise, leaders of governments have an imperfect assessment of each other's actions. To sustain cooperation, firms form trade associations and heads of state hold regular meetings to share information and coordinate their actions. How can efficient outcomes be sustained in an environment with imperfect monitoring? We address this question with a controlled experiment that assesses two canonical information management institutions: information delay and communication.

Delay of information is ubiquitous. Company bonuses to CEOs are given on a yearly basis. The G20 meetings, which from 2008 to 2011 were held on a semi-annual basis, now take place annually. The Kyoto Protocol, a global initiative to reduce emissions of greenhouse gases, establishes two commitment periods for the member countries: 2008-2012 and 2013-2020. A theoretical justification for delay is that it can help overcome the bounds on welfare imposed by inefficient provision of incentives. For concreteness, consider a repeated game with imperfect monitoring in which a noisy public signal of the chosen action profile arrives every period. The monitoring technology can be such that under public equilibria,<sup>32</sup> welfare is bounded away from efficiency by a substantial amount (Fudenberg et al., 1994; Sannikov and Skrzypacz, 2007). Near efficiency, however, is possible if the signal is delayed, i.e., if players receive several signals at a time instead of receiving a signal every period. This point was first made by Lehrer (1989) and Abreu et al. (1991) and has since become a standard technique in the theoretical literature on repeated games, especially in the study of both private monitoring and private strategies.<sup>33</sup> This literature exploits the delay of endogenous information: private signals and actions of other players. Intuitively, without delay, public equilibrium requires that both players are punished for a “bad” signal in every period it is observed. If the signal is shown every two

---

<sup>31</sup>Joint work with David Rahman.

<sup>32</sup>Intuitively, public equilibrium means behavior only depends on public information. By imperfect monitoring, past behavior is not public information, so players cannot react to their own past actions.

<sup>33</sup>See, for instance, Compte (1998); Ely et al. (2005); Hörner and Olszewski (2006); Kandori and Obara (2006).

periods, punishment can be triggered (with some probability) by *two* instances of bad news. Longer delays allow performance to be reviewed more efficiently.

Communication, likewise, is a pervasive element of human interaction, and experiments have shown that it leads to improved welfare and coordination (a review of this literature can be found in Crawford, 1998). Although most of the theoretical literature on communication or “cheap talk” (Crawford and Sobel, 1982) emphasizes limitations to information sharing when incentives are misaligned, the literature concludes that some communication is often better than none at all. In the context of repeated games, however, subjects’ payoffs are bounded away from efficiency if public equilibria are played, and this is true irrespective of the ability to communicate.<sup>34</sup>

We design the experiment around a prisoner’s dilemma game played repeatedly with imperfect monitoring and frequent actions elapsing at a rate of 0.15 seconds per period. Because one of our main goals is to see whether players take advantage of delay, this environment is particularly appropriate. With our chosen parameters, the efficient level of welfare is 30 and welfare levels above 20 cannot be sustained according in public equilibrium according to standard theory. In the treatments with delay, information arrives in 100 period blocks, making it possible to sustain welfare levels above 29. A game with a small number of periods would make the benefits of delay substantially less stark.

The experiment is described in detail in Section 3.3, but its basic features are the following. Subjects are randomly and anonymously matched into groups of two and earn points depending on the group’s chosen action profile. Instead of observing the other player’s actions, each subject observes a noisy public signal that has a positive drift if and only if both matched players cooperate. In the treatment without delay (treatment N), the public signal is shown in real time. In the treatment with delay (treatment D), the signal is shown in 100 period (15 second) windows. Two additional treatments allow subjects to their strategies with (treatment DC) and

---

<sup>34</sup>Some papers have explored communication as a useful tool for augmenting the set of equilibrium outcomes by allowing strategies to depend on the communicated information (e.g., Compte, 1998; Kandori and Matsushima, 1998; Obara, 2009). Kandori (2003) studies a repeated game with public monitoring and proves a folk theorem with communication in this environment. Although Kandori’s result requires more than two players, one can ask if a version of his solution is applicable to our example. It is shown in Rahman (2013a), however, that a folk theorem in public communication equilibria requires that the drift when both players defect differs from the drift with unilateral defection, a condition which is violated by our monitoring technology. The public equilibrium bound on payoffs persists.

without (treatment NC) delay of information.

With frequent actions, the bound on public Nash equilibrium payoffs can also be overcome with bounded rationality. If actions cannot be changed every period (e.g., because of physical constraints on reaction times), the players observe several signals before making their next decisions, and this bundling together of information makes welfare levels above 20 sustainable in equilibrium. For example, if it takes players five periods to respond, they can use a trigger strategy which starts off by cooperating and continues to do so as long as anything other than five bad signals is observed and defects with some probability if five bad signals are observed. To test whether bounded rationality affects behavior, we introduce a slow treatment (treatment S) that is identical to the baseline no delay, no communication case in all respects by two: a period lasts for a whole second, rather than 0.15 seconds, and the exchange rate between points and dollars is adjusted to equalize earnings per unit of time.

Our main results are the following:

**Result 1.** DELAY OF INFORMATION HINDERS COOPERATION.

**Result 2.** COMMUNICATION IMPROVES COOPERATION, ALLOWING PLAYERS TO EXCEED THE PUBLIC EQUILIBRIUM BOUND ON PAYOFFS.

**Result 3.** GIVING PLAYERS MORE TIME TO THINK ABOUT THEIR CHOICES HAS NO EFFECT ON BEHAVIOR.

The finding that delay leads to a decrease in welfare cannot be explained by (public)  $\varepsilon$ -equilibria, where each player is a small distance away from playing a best response to the other player's strategy.<sup>35</sup> Friedman and Oprea (2012), the first paper to systematically examine behavior in a continuous time prisoners' dilemma (with perfect monitoring), provides a useful reference point for this observation. The paper finds median cooperation rates above 90% in continuous time and provides a theoretical model to explain this data, building on earlier work by Radner (1986) and Simon and Stinchcombe (1989). Focusing on cut-off strategies  $K(s)$  with conditional cooperation until time  $s$  and unconditional defection thereafter, the authors show that  $\varepsilon$ -equilibria are consistent with their experiment's results.

Our results provide several counterpoints to this conclusion. First, when imperfect

---

<sup>35</sup>This is because the result that near efficiency can be sustained with delay is robust to small mistakes.

monitoring is introduced in an otherwise similar environment, (public)  $\varepsilon$ -equilibria cannot explain the observed behavioral regularities. The heart of [Friedman and Oprea](#)'s argument was that frequent actions permitted players to punish deviations quickly, rendering them unprofitable. With imperfect monitoring, it takes time to recognize a deviation, and as a result reacting quickly loses its power. In fact, we found no significant difference in cooperation rates between our slow (1 period per second) and fast (6 periods per second) treatment. Arguably, subjects cannot react to information at the rate of 1/6th of a second, so via this form bounded rationality, their ability to react promptly to deviations was limited more in the fast treatment than in the slow treatment. Nevertheless, cooperation rates were not significantly different. Secondly, since [Friedman and Oprea \(2012\)](#) studied perfect monitoring, deviations could be detected precisely in their experiments. On the other hand, in our experiments imperfect monitoring made it impossible for players to detect perfectly the behavior of their opponents. They needed repeated observations to make confident judgments regarding their opponents' behavior.

That delay of information leads to significant losses in welfare is our paper's main contribution. It has been pointed out that the efficiency gains associated with delay in [Abreu et al. \(1991\)](#) may in practice be counteracted by the benefits of receiving frequent feedback ([Levin, 2003](#)). In the context of a laboratory experiment, we find that this is indeed the case. Our results are in broad agreement with important findings in the industrial organization literature, which treats communication and information sharing as canonical ways of sustaining collusion ([Feuerstein, 2005](#)). In this line of research, there exist important examples of collusive institutions that choose *not* to delay noisy information. The Joint Executive Committee, a well-known railroad cartel which controlled much of railroad shipment in late nineteenth century United States, published weekly statistics that allowed cartel members to check on each other weekly ([Ulen, 1980](#)). Indeed, according to [Porter \(2005\)](#), "the cartel formation process [...] involves more than the issues studied in the repeated games literature. Dampening the short run incentives to cheat is only one facet of a cartels problems." We agree with this assessment and take our experimental results to point to the following basic fact: Contrary to standard theory, management of exogenous information can decrease welfare, while an institution that allows for additional information to be generated endogenously can lead to significant welfare benefits.

## 3.2 Related work

### 3.2.1 Experimental literature

There is a small but growing experimental literature on repeated games played with frequent actions. [Friedman and Oprea \(2012\)](#) showed that cooperation rates in a prisoner’s dilemma are higher when the game is played in quasi-continuous time than when time is discrete. [Bigoni et al. \(2011\)](#) compared the effects of fixed and random termination times in the same setting, extending related experiments of [Dal Bó \(2005\)](#) conducted in discrete time. [Oprea et al. \(2011\)](#) used a continuous “hawk-dove” game in an experimental test of evolutionary game theory. We follow the basic methodology established in these studies: An action is assumed to be fixed until changed by the subject, while payoff stocks are updated every period, which in our case lasts 0.15 seconds.

While subjects observed their partner’s choices in these studies, other experiments, in both discrete and continuous time, made use of imperfect monitoring. [Aoyagi and Fréchette \(2009\)](#) showed that welfare decreases in a repeated prisoner’s dilemma as the public signal becomes more noisy. [Ambrus and Greiner \(2012\)](#) studied the relationship between welfare and the severity of a punishment technology in a public good game. [Bigoni et al. \(2012\)](#) found that action frequency has a nonlinear impact on collusion when payoffs are updated in a quasi-continuous manner and monitoring is noisy.

Our study is the first to implement imperfect monitoring and information delay in a theoretically structured manner. [Cason and Khan \(1999\)](#) delayed the announcement of other subjects’ contributions in a public good game, interpreting information delay as an imperfect monitoring technology. As pointed out in [Aoyagi and Fréchette \(2009\)](#), such an interpretation of imperfect monitoring is at odds with the way the former is construed in theory. Moreover, information aggregation is irrelevant in a setting without noise.

Our experiment also manipulates subjects’ ability to communicate. The experimental literature on communication is vast and dates back to at least [Dawes et al. \(1977\)](#). Studies in this line typically find that communication increases cooperation rates amongst experimental participants. This finding, however, comes with some qualifications. [Charness \(2000\)](#), for instance, found that “minimalist” communication



protocols that allow players to announce their strategies are ineffective at improving cooperation rates in a one shot prisoner’s dilemma. Ben-Ner et al. (2007) found that numerical messages are much less effective than verbal ones at encouraging trusting and trustworthy behavior in a trust game. Charness et al. (2012) employed a design manipulating the subjects’ ability to communicate in a freeform manner and the rate at which periods elapsed, and found that communication had a much greater effect on contributions in continuous than in discrete time. This study relates to ours only loosely. First, it utilizes a setting with *perfect* monitoring. Second, the communication technology employed in our study is closer to the “minimalist” protocol of Charness (2000) or the numerical protocol of Ben-Ner et al. (2007) than the type of free-form communication employed in Charness et al. (2012).

### 3.2.2 Theoretical literature

The theoretical literature on repeated games with frequent actions is also small, recent and growing. With perfect monitoring, Simon and Stinchcombe (1989) developed an influential idea for sustaining cooperation in a Prisoners’ Dilemma with finite horizon, assuming a form of bounded rationality. Specifically, they assume that players can only react to their observations with some fixed delay. As actions become arbitrarily frequent, for any fixed delay in reaction times they show that it is possible to sustain cooperation in a Prisoners’ Dilemma—even if the horizon is fixed and finite.

Although Radner (1986) focused on the discrete time case, his results apply<sup>36</sup> just as much to games with frequent actions. He points out the discontinuity in the equilibrium payoff correspondence from an arbitrarily large but finite horizon to an infinite horizon, and then offers three different ways of restoring continuity. First, by introducing reputation, as in the famous “gang of four” papers (e.g., Kreps et al., 1982), cooperation becomes possible. Second, relaxing the behavioral predictions to  $\varepsilon$ -equilibria allows for some cooperation in equilibrium. Third, if players’ strategies are subject to being “executed” by finite state automata with a fixed upper bound on their number of states, then continuity of the equilibrium payoff correspondence is again restored with respect to the horizon.

Friedman and Oprea (2012) use these interesting results to understand their experimental results in a theoretically structured manner. They emphasize a combination

---

<sup>36</sup>Radner (1986) attributes some of the findings in his paper to others; see his paper for references.

of  $\varepsilon$ -equilibrium and delayed reaction as a way of explaining the behavior of subjects in their experiment. However, none of the arguments mentioned above generalize immediately to games with imperfect monitoring, and no such extension exists in the literature. Such a generalization is an interesting open problem that we leave for future research. On the other hand, our results seem to rule out both  $\varepsilon$ -equilibrium and finite automata arguments as drivers for cooperation in the Prisoners' Dilemma with imperfect monitoring: In theory, delay ought to add value even in  $\varepsilon$ -equilibrium and regardless of the feasible complexity of a strategy.

The existing literature on repeated games with imperfect monitoring has a long history by now, perhaps most notably Radner et al. (1986), Abreu et al. (1986, 1990), Abreu et al. (1991), as well as Fudenberg et al. (1994). Relatively recently, Sannikov (2007); Sannikov and Skrzypacz (2007, 2010) and Fudenberg and Levine (2007, 2009) extended these techniques and results to games with both frequent actions and imperfect monitoring. A crucial assumption that is made in all of these papers is that players behave according to (perfect) public Nash equilibrium. This restriction on the set of equilibria facilitates formal analysis of sustainable payoffs, often with stark behavioral predictions. For instance, according to Sannikov and Skrzypacz (2007), collusion is impossible in a repeated Cournot oligopoly with flexible production, and the amount of cooperation that is sustainable in the Prisoners' Dilemma is severely limited in public Nash equilibrium. This is shown in Proposition 1 below.

There is substantial theoretical precedent for the question of how exogenous delay of information helps players sustain cooperation. Starting from Lehrer (1989) and Abreu et al. (1991), the idea that players can attain better social outcomes by delaying and lumping information into blocks has been widely accepted and applied in various contexts. Thus, “block” strategies have been used to sustain socially desirable outcomes in Kandori and Matsushima (1998), Compte (1998), Obara (2008), Ely et al. (2005), Sugaya (2010) and others. Although none of these “cooperative institutions” survive in repeated games with frequent actions, it can be shown that even with frequent actions delay can still help (Rahman, 2013b)—at least in theory.

### 3.3 Experimental design

All treatments were programmed using zTree (Fischbacher, 2007b) and implemented with a between-group design following all standard practices of experimental eco-

nomics. The experiment had four treatments, described in detail below. Upon signing their consent forms, subjects in every treatment obtained a paper copy of the instructions and were shown a pre-recorded Power Point presentation explaining their task. They then played a two player repeated Prisoner’s Dilemma with imperfect monitoring and frequent actions through a computerized interface. Appendix 4.9 provides the instructions to the NC treatment.<sup>37</sup>

In all treatments other than treatment S, a time period lasted  $\Delta t = 0.15$  seconds. At the beginning of each match, subjects chose between pressing an orange button (“cooperate”) and a purple button (“defect”). After their initial choice, they could change their selection at any time and as often as they wanted. Not pressing any buttons during a time period amounted to maintaining their last recorded choice. Following one practice match, subjects were randomly and anonymously matched several times. In every period, the probability of a match terminating was  $1/700$ . Following Murnighan and Roth (1983), we identify the continuation probability with the discount factor. The first match to end after 45 minutes elapsed since the beginning of the experiment marked the end of a session. If subject were in mid-match at 50 minutes after the experiment begin, we overrode the random termination rule and terminated the match randomly. Payment consisted of the final payoff from a randomly selected match, converted from points to dollars at the exchange rate of 40 (treatments N and D) or 20 (treatments NC and DC) points per cent.

Depending on whether she chose to cooperate ( $C_{it}$ ) or defect ( $D_{it}$ ), subject  $i$ ’s stock of points increased by  $u_{it}$  in period  $t$  according to the following table:

	$C$	$D$
$C$	15, 15	0, 20
$D$	20, 0	2, 2

Subjects did not find out their earnings until the end of the session, at which point they received their accumulated earnings  $\sum u_{it}(a_t)$  from every match. Instead of observing her partner’s actions, each subject was shown a public signal that could go up with probability  $p(a_t)$  or down with probability  $1 - p(a_t)$ , with  $p(a_t)$  determined as in the table below:

<sup>37</sup>The instructions to treatments with delay differ in that the sentence “The process will be displayed in real time, in blocks of 100 periods” in the “Information” section is replaced by “You will only observe the evolution of this process at the end of each block of 100 periods.” The instructions to treatments N and D omit the “Communication” section.

	<i>C</i>	<i>D</i>
<i>C</i>	$\frac{3}{4}$	$\frac{1}{2}$
<i>D</i>	$\frac{1}{2}$	$\frac{1}{2}$

Conditional probability  $p$  of the good signal

In treatment N, information arrived continuously, and in treatment D it arrived at the end of each 100 period (15 second) block, when it was revealed in five three second-long lumps. In both treatments, points were converted to dollars at an exchange rate of 40 points per cent. In treatment S, information also arrived continuously, but a time period lasted for  $\Delta t = 1$  second. The continuation probability (discount factor) was identical to that in the other treatments, but points were converted to dollars at an exchange rate of 6 points per cent. This ensured that earnings per unit of time were the same.<sup>38</sup> The slow treatment also had 150 (unpaid) periods in the practice match, compared to 250 in the other treatments. This ensured that the practice match does not go on for an unnecessarily long length of time, but that the players still get an opportunity to experience a match with more than one 100-period block.

Treatments NC and DC allowed subjects to communicate cut-off strategies in an environment without (NC) and with (DC) delay of the public signal. Because communication made each match longer, we introduced a more generous exchange rate in these treatments to help smooth out earnings across sessions: 20, rather than 40 points, were converted into a cent. In all treatments of the experiment, subjects pressed a “continue” button every 100 periods—at the end of each 15 second block. This block structure was introduced to minimize the difference between the differences between ways in which information is presented in treatments with and without delay. Note that it leaves our theoretical predictions unaltered.

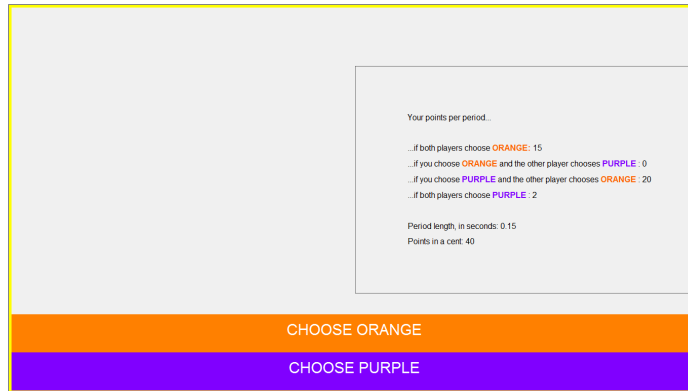
The communication in these treatments took place at the beginning of each match, before subjects took their initial actions, and at the end of each 100 period block. At the beginning of the match, subjects provided answers to the following questions:

- This block, I will choose ORANGE this percentage of the time: \_\_\_%
  
- If this block’s signal position is [above/below] the number \_\_\_, I will respond

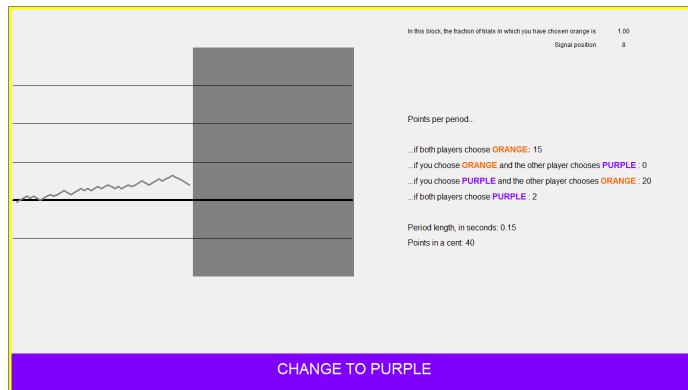
---

<sup>38</sup> $(1/40)$  cents/point x 15 points/period x  $(1/.15)$  periods/second = 2.5 cents/second. To make the slow world the same, change the exchange rate to  $((1/40) \times (1/.15))$  cents/point x 15 points/period x 1 period/second = 2.5 cents/second.

Beginning of match:



Mid-block:



End of block:

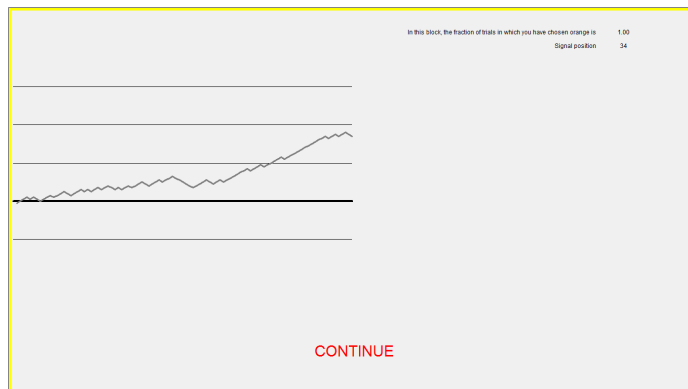


Figure 5: Screenshots of the first block in treatment N (no delay, no communication). In the beginning of each match, a subject selects her initial action (top). She then plays a prisoner's dilemma game with imperfect monitoring, in which a common signal of the chosen action profile is observed (middle, bottom). The player is allowed to change her action as often as she desires in the course of the game (middle). Every 15 seconds, the player presses a "continue" button (bottom), which erases past signals from the display.

by choosing ORANGE this percentage of the time in the following block: \_\_\_%.  
Otherwise, I will choose ORANGE this percentage of the time: \_\_\_%

After everyone submitted their answers and initial actions, subjects saw their partner's answers displayed on the screen for 30 seconds. When this screen timed out, the game began to elapse. At the end of each 100 period block, subjects were given an opportunity to revise their answers. After everyone's new answers were submitted, each subject looked their partner's new answers for 15 seconds before the next block started.

### 3.4 Theoretical predictions

From the point of view of the theory, the modeling choice of imperfect monitoring with frequent actions is useful for three reasons. First, it is consistent with many real-world applications. Second, it disciplines the design of institutions considerably by forcing them to be robust to the friction of a fixed period length by discouraging infinitesimal deviations.<sup>39</sup> Third, it delivers a mathematically tractable analysis of the problem. We now describe the theoretical considerations relevant to our study.

#### 3.4.1 Public equilibrium payoffs

The study of public equilibria is practically the norm in repeated games, especially those with frequent actions. As such, it is important to understand the restriction that such equilibria impose on equilibrium payoffs. Fortunately, their recursive nature deliver a simple, partial identification for public equilibria: the maximal payoff under public equilibria is given by 20 points per period. This claim is proved in the following proposition, see [Figure 6](#) below for a graphical illustration.

**Proposition 1.** *Let  $\gamma(\delta) = \max\{v_1 + v_2 : v \in E(\delta)\}$ , where  $E(\delta)$  is the set of public equilibrium payoff vectors of the game with discount factor  $\delta < 1$ . Then,  $\gamma(\delta) \leq 20$  for every  $\delta$ . This bound continues to hold in public communication equilibrium.*

---

<sup>39</sup>This rules out most of the institutions in the game theory literature, including [Abreu et al. \(1991\)](#) to [Compte \(1998\)](#), [Kandori and Matsushima \(1998\)](#), [Ely et al. \(2005\)](#), [Hörner and Olszewski \(2006\)](#), [Sugaya \(2010\)](#) and beyond.

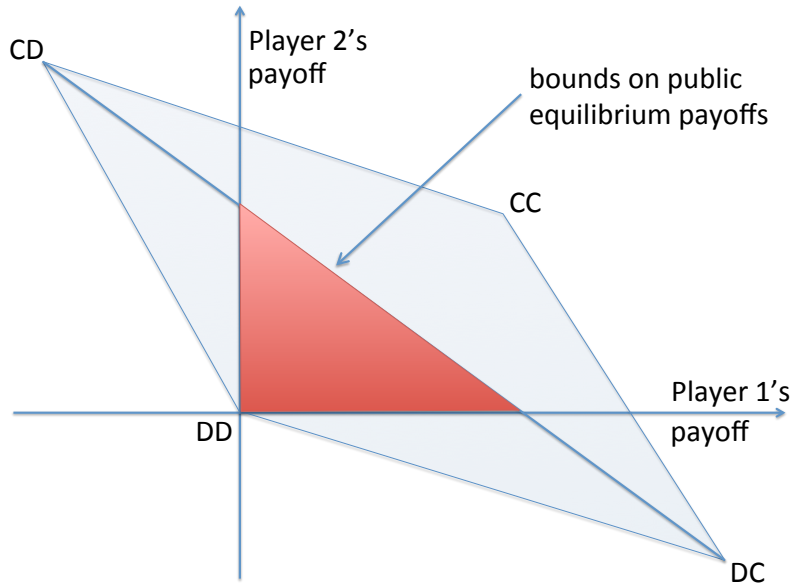


Figure 6: Flow payoffs in the Prisoners' Dilemma with two players

Note that [Proposition 1](#) delivers the same theoretical bound for both public Nash equilibrium and public communication equilibrium. This is due to the particular choice of information structure in our experiment, where the probability of good news is the same if there is only one cooperator or none at all. Therefore, the fact that our treatment with communication exceeded the welfare bound of 20 is not consistent with public communication equilibrium.

### 3.4.2 How information delay can help

The insight that lumping information together may improve incentives is not new, dating back at least to [Lehrer \(1989\)](#).<sup>40</sup> For our purposes, the construction due to [Abreu et al. \(1991\)](#) is a particularly useful way of describing it.

Suppose that, instead of the signal arriving every period, it was possible to lump the information in such a way that the signal only arrived at the end of every  $T$ -period block. [Abreu et al. \(1991\)](#) show how players can improve upon a welfare of 20 by delaying information this way. Consider the following strongly symmetric strategies, to be called *AMP block strategies*. Every player cooperates for  $T$  periods. At the end of the  $T$ -period block, the  $T$  public signals for each period in the block

<sup>40</sup>[Lehrer \(1989\)](#) studied repeated games without discounting, though.

arrive to the players. If every signal was bad then continuation play consists of mutual defection henceforth with some probability  $\alpha$ . Otherwise, they continue to cooperate for the next block with the same contingency.

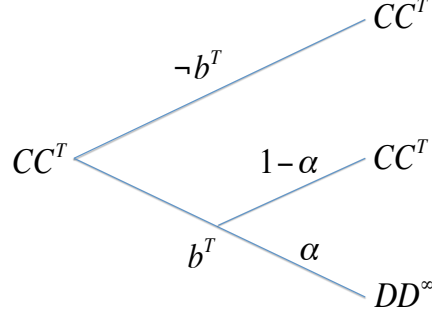


Figure 7: AMP block strategies ( $b^T = T$  bad signals)

The probability of  $T$  consecutive bad signals equals  $q_2^T$  in equilibrium, that is, assuming mutual cooperation throughout the block. A player's lifetime utility under this strategy profile is therefore given by

$$v = (1 - \delta^T)15 + \delta^T [(1 - q_2^T)v + q_2^T((1 - \alpha)v + 2\alpha)].$$

Rearranging,

$$v = 15 - \frac{\delta^T}{1 - \delta^T} q_2^T \alpha (v - 2). \quad (3.1)$$

Discouraging a deviation in the very first period of the block requires that the utility gained from defecting,  $(1 - \delta)5$ , be outweighed by the associated loss in continuation payoff. This is given by the change in probability of punishment from the one-period deviation,  $q_2^{T-1}(q_1 - q_2)$ , times the opportunity cost of punishment,  $\delta^T \alpha (v - 2)$ . Since  $q_1 - q_2 = .25 = q_2$ , this incentive constraint may be written as

$$(1 - \delta)5 \leq \delta^T q_2^T \alpha (v - 2). \quad (3.2)$$

A key insight behind the welfare properties of AMP block strategies is that discouraging one deviation discourages all others, as the next result shows. The intuition for it is this. The gains from deviating grow linearly, whereas the costs grow exponentially in the number of deviations. Therefore discouraging one deviation discourages them all.



**Lemma 1.** *If the AMP block strategies above discourage a deviation in any single period of a block then they discourage every deviation, that is, they constitute an equilibrium.*

Consider maximizing  $v$ , the strongly symmetric equilibrium payoffs above, with respect to  $\alpha$  such that the AMP block strategies above remain an equilibrium. At an optimum, the incentive constraint (3.2) must bind, since otherwise by (3.1) we would be able to feasibly lower  $\alpha$  further and increase  $v$ , contradicting optimality. If (3.2) binds then the maximum value of  $v$  equals

$$v = 15 - 5 \frac{1 - \delta}{1 - \delta^T}.$$

On the other hand, feasibility requires that  $\alpha \leq 1$ , since it is a probability. Substituting for  $v$  and this inequality in (3.2) and rearranging gives

$$5(1 - \delta) \left[ \frac{1}{1 - \delta^T} + \frac{1}{(\delta q_2)^T} \right] \leq 13. \quad (3.3)$$

This inequality places a restriction on the exogenous parameters of the game for the strategy profile above to be an equilibrium. [Abreu et al. \(1991\)](#) used a version of this bound to argue a result along the following lines.

**Proposition 2.** *For every block length  $T \in \mathbb{N}$ , there exists  $\delta < 1$  sufficiently large that the strategies above constitute an equilibrium. Moreover,*

$$\lim_{T \rightarrow \infty} \lim_{\delta \rightarrow 1} v = 15.$$

Note that  $15 + 15 > 20$ . This shows that the public equilibrium bound on payoffs can be overcome with delay.

### 3.4.3 Delay with practical cut-offs

AMP equilibria suffer from the following basic problem: as  $T$  grows, the demands placed on  $\delta$  for an equilibrium become unreasonable. For instance, if  $T = 20$  (or 3 seconds) then (3.3) requires  $\delta$  to be so unreasonably close to 1 that interpreting it as the probability of termination would imply an expected duration of a match to be more than 1,000 years, using the parameters from our experiment.

Essentially, the construction by [Abreu et al. \(1991\)](#) is too lenient over when players are punished. It takes  $T$  bad signal realizations to trigger punishment. A less lenient, but more practical, approach to making use of delay is the following one taken from [Rahman \(2013b\)](#). If  $T$  is large then the distribution of signals is close to normally distributed by the Central Limit Theorem, and the likelihood that in equilibrium there will be more than  $\frac{1}{2}T$  bad signals is relatively low, since the expected number of bad signals is  $\frac{1}{4}T$ .<sup>41</sup> Therefore, punishment actually occurs relatively infrequently in equilibrium. Nevertheless, it is still the case that discouraging a single deviation discourages them all, which helps to establish a Folk Theorem.

Specifically, consider the following strategies, called *practical cut-off strategies*. As with AMP block strategies in [Section 3.4.2](#), players cooperate over a  $T$ -period block with information delay. If the number of bad signals at the end of the block is greater than  $\frac{1}{2}T$ , then players choose mutual defection henceforth with some probability  $\alpha$ . Otherwise, they continue to cooperate in the next block, with the same contingent plan. The only difference between AMP block strategies and practical cut-off strategies is in the number of bad signals that trigger punishment. In the former, this number is  $T$ , whereas in the latter it is  $\frac{1}{2}T$ . Nevertheless, [Rahman \(2013b\)](#) shows that these cut-offs lead to efficient outcomes with reasonable discount rates, in contrast with the previous objection to AMP block strategies. This shows that—at least in theory—it is possible for players to substantially improve on the public equilibrium benchmark of [Proposition 1](#).

### 3.4.4 How bounded rationality can help

As noted in the introduction, bounded rationality can also help players overcome the bound on public equilibrium payoffs. Let  $\tau = 2$  denote the number of periods that a player is unable to change her action, and consider the following strongly symmetric strategies. Every player cooperates for  $\tau$  periods. At the end of the  $\tau$ -period block, the players consider cooperating if anything other than  $\tau$  bad signals is observed. If  $\tau$  bad signals are observed, the players switch to defection with some probability  $\alpha$ .

The probability of 2 consecutive bad signals equals  $q_2^2$  in equilibrium. A player's

---

<sup>41</sup>This cut-off is lower than might be expected from a version of the mechanism proposed by [Kandori and Matsushima \(1998\)](#) in discrete time, closer to  $\frac{1}{4}T$ , which would not approximate full cooperation due to a non-vanishing punishment probability.

lifetime utility under this strategy profile is therefore given by

$$v = (1 - \delta^2)15 + \delta^2[(1 - q_2^2)v + q_2^2((1 - \alpha)v + 2\alpha)].$$

Rearranging,

$$v = 15 - \frac{\delta^2}{1 - \delta^2} q_2^2 \alpha (v - 2). \quad (3.4)$$

Discouraging a deviation requires that the utility gained from defecting,  $(1 - \delta^2)5$ , be outweighed by the associated loss in continuation payoff. Thus, the incentive constraint is

$$(1 - \delta^2)5 \leq q_2^2 \alpha (q_1^2 - q_2^2)(v - 2). \quad (3.5)$$

Consider maximizing  $v$  with respect to  $\alpha$  such that the trigger strategies above remain an equilibrium. At an optimum, the incentive constraint (3.5) must bind, since otherwise by (3.4) we would be able to feasibly lower  $\alpha$  further and increase  $v$ , contradicting optimality. If (3.5) binds then the maximum value of  $v$  equals

$$v = 15 - \frac{5}{\left(\frac{q_1}{q_2}\right)^2 - 1} = 15 - 5/3 \approx 13.33 > 10.$$

It is easy to check that the feasibility constraint, that is,  $0 \leq \alpha \leq 1$ , is satisfied given the parameters of the experiment.

### 3.5 Results

Data was collected from 248 University of Minnesota undergraduate students at the Anderson Hall Social and Behavioral Sciences Laboratory. Table 13 reports select summary information.<sup>42</sup> To estimate the effect of our treatments on cooperation, we regressed each subject's average cooperation rate in non-practice matches<sup>43</sup> on three dummies: Delay (=1 for treatments with delay), Communication (=1 for treatments with communication) and Slow (=1 for treatment S). The regression was performed with session-clustered standard errors. The results, presented in the left column of

---

<sup>42</sup>Notice that treatments N and D had more matches than treatments NC and DC. This is because communication took up a significant portion of time in the latter treatments, as discussed in Section 3.3.

<sup>43</sup>The analysis here is restricted to non-practice matches. We discuss behavior in practice matches and the dynamics of cooperation below.

Table 14, show that delay of information led to lower and communication to higher cooperation rates.

		<i>Delay</i>		<i>No delay</i> <i>No communication</i> $\Delta t = 1$ sec.
		$\Delta t=0.15$ sec.	$c=15$ sec.	
<i>Communication</i>	No	Treatment N	Treatment D	Treatment S
		6 sessions	3 sessions	
	$N = 64$	$N = 46$		
	287 matches	169 matches		
Yes	Yes	Treatment NC	Treatment DC	4 sessions $N=42$ 49 matches
		3 sessions	3 sessions	
		$N = 44$	$N = 52$	
		109 matches	83 matches	

Table 13: Summary statistics of the experimental treatments

While other studies found that communication improves cooperation in games, our result that delay hinders cooperation is entirely new; it is therefore important to verify its replicability. We found a similar result in a pilot experiment that we conducted in the summer of 2013 with a sample size of 66 subjects, where under different parameters delay has a significantly negative effect on welfare ( $P < 0.05$  with session-clustered errors). This experiment also had frequent actions and imperfect monitoring, but subjects observed their own stocks of payoffs, rather than the noisy signal. In one treatment, the payoff stock was observed in real time, and the second treatment, it was observed with delay. The payoff stock, however, depended on one's own chosen action, the chosen action of the opponent, and noise. We highlight our first result below:

**Result 1.** DELAY OF INFORMATION HINDERS COOPERATION.

The effects of the treatment variables on subjects' average attained payoffs are shown in the right column of Table 14. It is apparent the effects on welfare exactly parallel the effects on cooperation rates described above. The average payoff in treatments N and DC does not differ significantly from 10 ( $P$ -values of 0.742 and 0.937, respectively); the average payoff in treatment D is significantly below 10 ( $P < 0.01$ ), and the average payoff in treatment NC is significantly above 10 ( $P < 0.01$ ). This latter result suggests off-equilibrium behavior, private strategies, or equilibrium behavior with bounded rationality in treatment NC. We highlight this result below:

**Result 2.** COMMUNICATION IMPROVES COOPERATION, ALLOWING PLAYERS

	Cooperation	Average payoff
Delay	-0.102*** (0.0359)	-1.251*** (0.440)
Communication	0.111*** (0.0364)	1.405*** (0.444)
Slow	-0.0264 (0.0392)	-0.315 (0.463)
Constant	0.558**** (0.0288)	9.882**** (0.353)
Observations	248	248

Session-clustered standard errors in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ , \*\*\*\*  $p < 0.001$

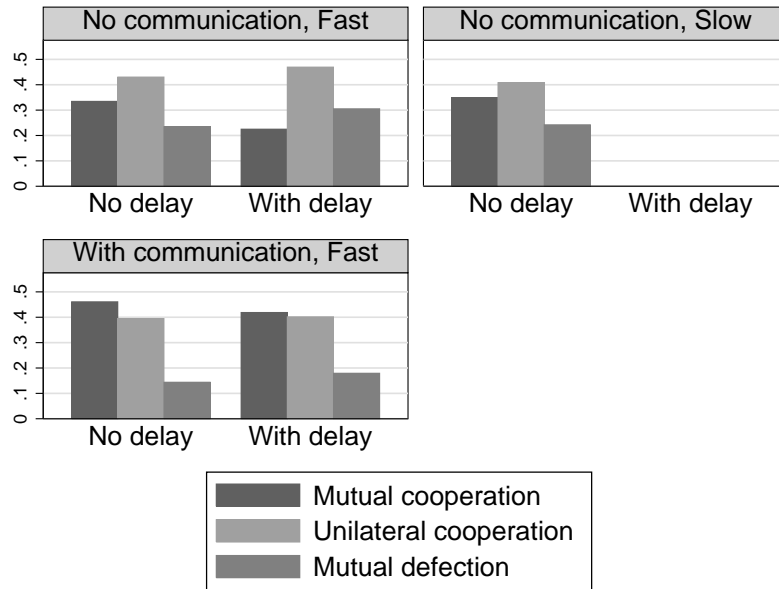
Table 14: Treatment effects on subjects' average cooperation rates and payoffs.

#### TO EXCEED THE PUBLIC EQUILIBRIUM BOUND ON PAYOFFS.

As noted previously, subjects' inability to change their action in every period of the treatments with frequent actions raises the question of whether this sort of bounded rationality has an effect on behavior. That the subjects exceed the bound on public Nash equilibrium payoffs in treatment NC makes this question all the more salient. We find, however, that slowing down the experiment has no significant effect on subjects' cooperation rates or payoffs (Table 14). In magnitude, the effect of the Slow dummy on individual cooperation rates is less than 3%, and the attained payoffs levels are virtually the same. This is our third significant result:

**Result 3.** GIVING PLAYERS MORE TIME TO THINK ABOUT THEIR CHOICES HAS NO EFFECT ON BEHAVIOR.

We also explored the effect of our treatment variables on *mutual* cooperation rates. To this end, we created three variables for each of the 697 matches in the dataset: percentage of time spent in action profile (C,C), percentage of time spent in action profile (C,D) or (D,C), and percentage of time spent in action profile (D,D). These profiles of cooperation rates are plotted for different treatments in Figure 8. In the bottom part of the figure, we report the results of regressions in which these variables



	(C,C)	(C,D) or (D,C)	(D,D)
Delay	-0.0889** (0.0419)	0.0293 (0.0271)	0.0596* (0.0329)
Communication	0.154*** (0.0491)	-0.0488 (0.0309)	-0.105*** (0.0280)
Slow	0.0224 (0.0527)	-0.0257 (0.0323)	0.00331 (0.0242)
Constant	0.327**** (0.0378)	0.434**** (0.0210)	0.239**** (0.0203)
Observations	697	697	697

Session-clustered standard errors in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ , \*\*\*\*  $p < 0.001$

Figure 8: Treatment effects on subjects' mutual cooperation rates.

are regressed against the treatment dummies. Delay made mutual cooperation less likely ( $P < 0.05$ ) and mutual defection more likely ( $P < 0.1$ ). Communication had the reverse effect, increasing mutual cooperation ( $P < 0.01$ ) and decreasing mutual defection ( $P < 0.01$ ). Interestingly, neither of these variables had an effect on unilateral cooperation rates. In the slow treatment, the profile of cooperation did not significantly differ from that in the baseline treatment (N), providing additional evidence for Result 3 above.

### 3.5.1 Dynamics

To study the dynamics of cooperation in the experiment, we changed the unit of analysis to average cooperation rate per block and re-ran the regression reported in Table 14 for all blocks and matches, including practice matches. We included a practice dummy, a block number variable (1 for the first 100 periods, 2 for the next 100 periods, etc.), and a match number variable. The results of this regression and an analogous regression in which the dependent variable is the subject's average payoff per block are reported in the left-most two columns of Table 15. Subjects cooperated significantly more often in the practice matches ( $P < 0.05$ ), but the match and block dummies failed to reach significance. The average cooperation rate in practice matches was approximately 65%, while the average in non-practice matches was approximately 55%.

We also re-ran the regressions reported in Table 14 for the practice matches. The results of these regressions are shown in the right-most two columns of Table 15. Unlike Table 14, Table 15 shows none of the treatment dummies as significant. I.e., subjects sustained comparably high cooperation rates and welfare levels in every treatment of the practice matches. Even in the slow treatment, where average payoffs above 10 cannot be sustained in public equilibrium with or without response time constraints,<sup>44</sup> we find that the average attained welfare level is significantly higher than 10 (mean=10.87,  $P < 0.05$ , session-robust standard errors). We summarize the findings described above as follows:

**Result 4.** THERE ARE SIGNIFICANT DIFFERENCES IN BEHAVIOR BETWEEN PRACTICE AND PAID MATCHES (AND EVIDENCE OF OFF-EQUILIBRIUM BEHAVIOR

---

<sup>44</sup>In this treatment, periods proceed at a slow enough rate that subjects can physically respond in every period.

	All data		Practice data	
	Cooperation (per block)	Average payoff (per block)	Cooperation (per subject)	Average payoff (per subject)
Delay	-0.0861** (0.0312)	-2.092** (0.777)	-0.0373 (0.0538)	-0.409 (0.651)
Communication	0.0936*** (0.0323)	2.302*** (0.785)	0.0590 (0.0542)	0.583 (0.655)
Slow	-0.0409 (0.0346)	-1.021 (0.815)	-0.00979 (0.0353)	-0.203 (0.443)
Practice	0.0634** (0.0260)	1.554** (0.661)		
Match	-0.00548 (0.00352)	-0.140 (0.0866)		
Block	-0.000817 (0.00118)	-0.0267 (0.0291)		
Constant	0.595**** (0.0300)	20.76**** (0.743)	0.648**** (0.0315)	11.08**** (0.386)
Observations	9674	4837	248	248

Session-clustered standard errors in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ , \*\*\*\*  $p < 0.001$

Table 15: Dynamics of cooperation in the experiment. There was significantly more cooperation in practice than in paid matches, but little learning once the paid matches start (left column). In the practice matches, none of the experimental treatments affected cooperation rates (right column).



IN PRACTICE MATCHES), BUT LITTLE LEARNING ONCE THE PAID MATCHES START.

To explore the dynamics of cooperation further, we examined whether what happens in the very first period of the match explains what happens in the rest of the match. To this end, we controlled for the subjects' chosen profile of cooperation in the regressions reported in Figure 8 by introducing two dummy variables: a mutual cooperation dummy (=1 if the players chose (C,C) in the first period of the match) and a unilateral cooperation dummy (=1 if the players chosen (C,D) or (D,C) in the first period of the match). The results are shown in Table 16. Both of the new variables have highly significant coefficients. Moreover, the effect of delay loses significance both when mutual cooperation in the match ( $P = 0.202$ ) and mutual defection in the match ( $P = 0.595$ ) are dependent variables. The effect of communication remains significant but falls in magnitude, from 15.4% to 9.13% when the percentage of match spent in mutual cooperation is the dependent variable, and from -10.5% to -6.53% when the dependent variable is mutual defection. Thus, first period profiles of cooperation largely explain the effects of the treatment variables. We highlight this result below:

**Result 5.** THE EFFECT OF DELAY (COMMUNICATION) ON COOPERATION IN THE COURSE ON THE MATCH IS LARGELY (PARTIALLY) EXPLAINED BY DIFFERENCES IN FIRST PERIOD PROFILES OF COOPERATION.

### 3.5.2 Strategies

We next looked at how subjects' behavior depended on the realized public signals. We averaged actions taken in each period (1-100) across all sessions, matches and blocks, taking into account the evolution of the public signal in the previous block. The average cooperation rates following blocks with different numbers of good news are plotted in Figure 9. This figure makes apparent a number of behavioral regularities.

First, it suggests that subjects cooperate more in blocks that follow blocks with high realizations of the public signal, and that this is true for treatments with and without delay. The figure obscures the effect of the public signal, however, because past realizations of the signal are related to past cooperation rates. We therefore ran regressions in which a subject's cooperation rate in block  $b$  was regressed against the public signal at the end of block  $b - 1$ , using the partner's cooperation rate in block

	(C,C)	(C,D) or (D,C)	(D,D)
Delay	-0.0435 (0.0330)	0.0281 (0.0226)	0.0154 (0.0286)
Communication	0.0913** (0.0336)	-0.0260 (0.0244)	-0.0653*** (0.0213)
Slow	0.0325 (0.0800)	-0.0381 (0.0486)	0.00564 (0.0355)
(C,C) in first period	0.414**** (0.0380)	0.0396 (0.0400)	-0.453**** (0.0513)
(C,D) or (D,C) in first period	0.0761*** (0.0218)	0.255**** (0.0393)	-0.331**** (0.0486)
Constant	0.124*** (0.0377)	0.297**** (0.0356)	0.579**** (0.0488)
Observations	697	697	697

Session-clustered standard errors in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ , \*\*\*\*  $p < 0.001$

Table 16: The action profile chosen by the players in the beginning of the match had a significant effect on the average action profile in the match. Note that the effect of the delay variable loses significance when players' first period choices are controlled for.

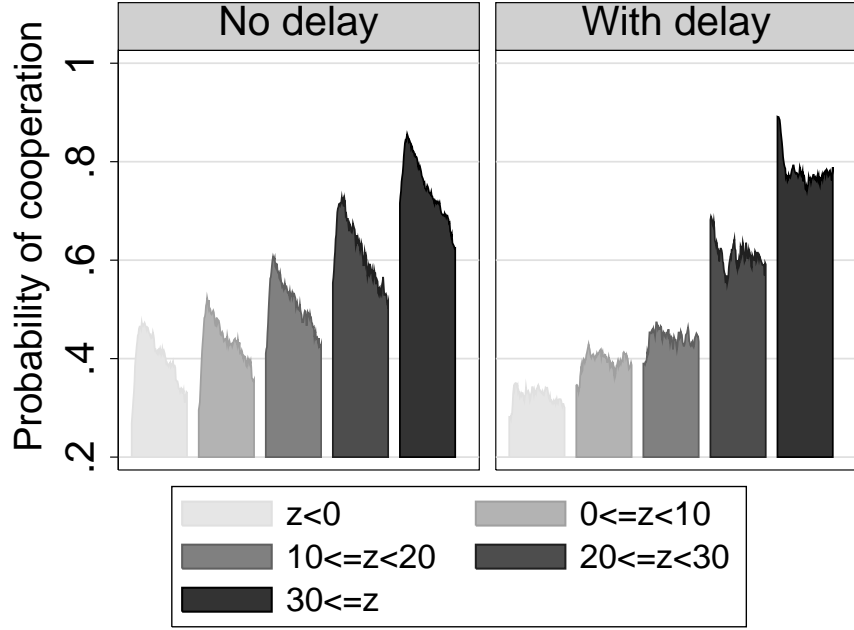


Figure 9: Behavior and the public signal. Each bar in plots average cooperation in rates in periods 1-100 of the block, with period number on the horizontal axis. Colors code the variable  $z$ , defined as the number of good news less the number of bad news received in the previous block.

$b - 1$  as in instrument. Specifically, the model we estimated was

$$a_{ib} = \beta z_{i,b-1} + \alpha_i + \epsilon_{ib},$$

where  $a_{ib}$  is subject  $i$ 's average cooperation rate in block  $b$  and  $z_{i,b-1}$  is the number of good news minus the number of bad news observed in the previous block  $b - 1$ . Because  $\epsilon_{ib}$  may be related to  $a_{i,b-1}$ , which is strongly correlated with  $z_{i,b-1}$ , there is a potential endogeneity issue. To address it, we used  $a_{-i,b-1}$ , the cooperation rate of  $i$ 's partner in block  $b - 1$ , as an instrument for  $z_{i,b-1}$ . Because  $i$  never observes her partner's actions except through the public signal, the instrument is valid.

We ran these regressions separately for every treatment, including subject fixed effects as covariates, and clustering the standard errors by session. The results are shown in Table 17. Once past cooperation rates are controlled for, we find that past realizations of the signal have highly significant positive effects on behavior ( $P < 0.001$ ) in every treatment without communication, and in the treatment with communication and no delay. The positive and highly significant relationship between signals and behavior constitutes our sixth major finding:

Treatment	N	D	S
Prev. block signals	0.00477****	0.00530****	0.00154****
Prev. block signals	(0.000631)	(0.000590)	(0.000388)
Observations	2994	1776	362

Treatment	NC	DC
	0.00316****	0.00114
	(0.000756)	(0.00249)
Observations	676	810

Session-clustered standard errors in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ , \*\*\*\*  $p < 0.001$

Table 17: Behavior and signals received in the previous 100-period block. The partner’s cooperation rate in the previous block is used as an instrument for previous block’s final signal position. The first stage F-statistics in these regressions range from 151.346 to 2389.576.

**Result 6.** SUBJECTS PROVIDE EACH OTHER WITH INCENTIVES BY COOPERATING MORE AFTER RECEIVING GOOD NEWS.

Note that the coefficient on the signal variable is almost two times smaller in magnitude in the treatments with communication. This suggests that when the subjects are able to communicate, they rely on old signals less, and that the messages have informational content. To confirm this, we looked at correlations between what the subjects communicated and how they behaved. For the 96 subjects who participated in the treatments with communication, the average promised cooperation rate in the non-practice matches was approximately 63% (promises averaged across blocks for each subject). The actual cooperation rate, in comparison, was approximately 69%. The correlation between promised and actual cooperation rates is strong and highly significant (correlation coefficient of .6214,  $P < 0.001$ ). We interpret this as strong evidence that messages in our data relate to behavior:

**Result 7.** THE MESSAGES HAVE INFORMATIONAL CONTENT.

### 3.5.3 Periodicity of behavior

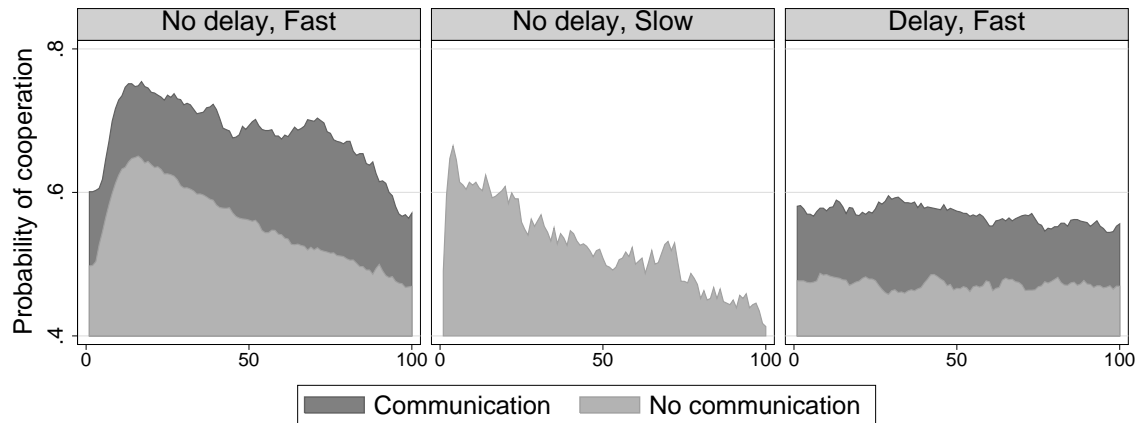
Figure 9 also suggests that behavior in treatments *without delay* has a highly periodic nature: regardless of the number of good signals observed in the previous block, subjects in these treatments are more likely to defect in the late periods of each block than they are in the early ones. This can also be seen in Figure 10, which aggregates the data for different past realizations of the public signal and plots the per-period cooperation rates for different treatments. Cooperation in these treatments is periodic with and without communication (top panel), with no clear time trend across blocks (bottom panel). Remarkably, cooperation rates in the slow treatment increase at about the same rate (per period) as they do in the fast treatments without delay, and then follow a similar gradual decline.

**Result 8.** BLOCKS IN TREATMENTS WITHOUT DELAY EXHIBIT A STRIKING PERIODICITY IN BEHAVIOR.

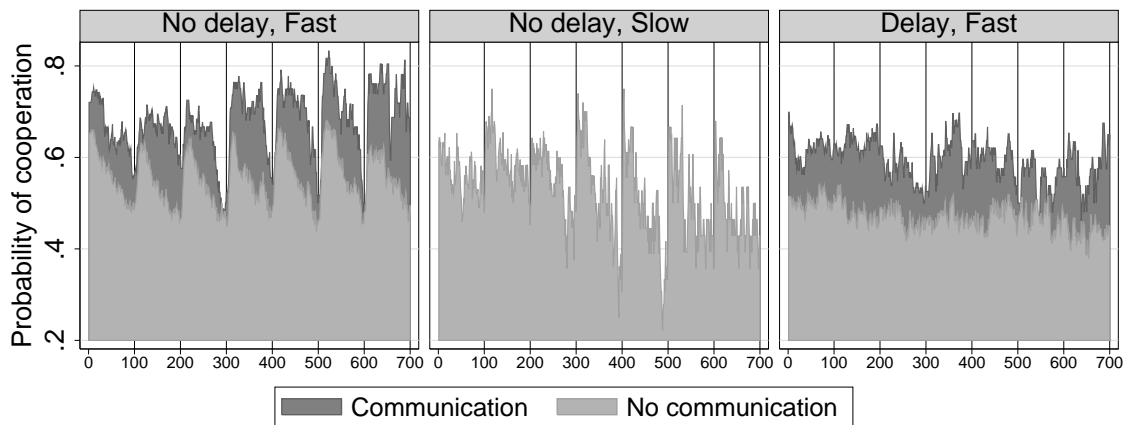
Note that the slope of the within-block trend is different in treatments with and without communication. Without communication, cooperation declines at a constant slope. With communication, it stays at a high level for a longer period of time, and then declines steeply. This observation sheds light on some of the dynamics of cooperation described above. Recall that the first period cooperation profile can explain most of the treatment effect of delay, but that the effect of communication remains significant even when the first period cooperation profile is controlled for. That cooperation declines faster without communication provides one avenue for a within block (and hence within-match) effect of the communication treatment.

One interpretation of this observed periodicity in behavior is that it reflects a focal point introduced by our design, which divides the match into 100-period blocks. This focal point may influence subjects to coordinate on mutual cooperation around the beginning of each block. Intuitively, a break (like a postman knocking on the door of a couple that is fighting) distracts the subjects, thereby restoring initial levels of cooperation. However, because behavior is not periodic in treatments with delay, this explanation is problematic.

An arguably more plausible interpretation is that players observe signals and choose their level of cooperation on the basis of the progress of the signal. For instance, suppose that in equilibrium the drift of the signal is fixed at some level while players are in a cooperative phase. As soon as the Brownian motion passes through a



(a) The average block



(b) First seven blocks of the average match

Figure 10: Periodicity of behavior. The top panels plot data averaged across all blocks and matches. In treatments without delay, cooperation rates start out high, decline over time, and refresh at the beginning of the next 100 period block. The decline in cooperation is faster for treatments without communication. No periodicity of behavior is observed in the treatments with delay. Behavior in the slow treatment follows the same pattern as behavior in the treatment with frequent actions and no communication.

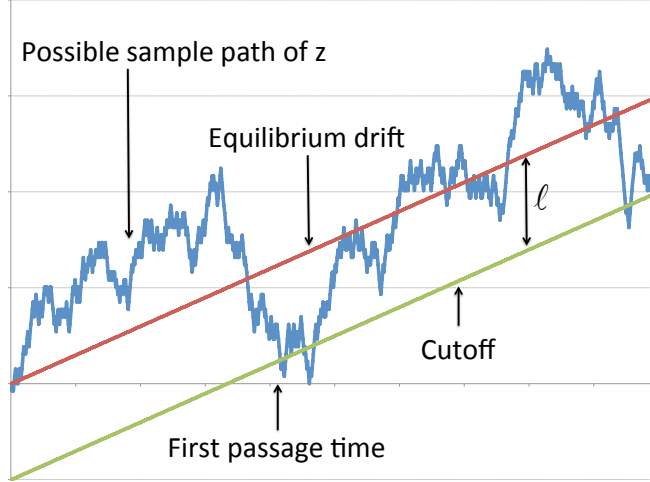


Figure 11: Illustration of strategies with first passage times determining switches in cooperation rates. The parameter  $\ell$  describes the amount of leeway a player gives to the public signal before switching from a cooperative regime to a more defective regime.

window of some length  $\ell$  from this drift, players move to a punishment phase for the rest of the block, and then restart their strategy in the next block, perhaps updating  $\ell$  on the basis of the length of time it took for the Brownian motion to escape the previous cutoff in the last period. See Figure 11. Assuming a common threshold  $\ell$  across matches and subjects, the density  $f$  of first passage times would take the form

$$f(t) = \frac{\ell}{\sqrt{2\pi t^3}} \exp\left(-\frac{\ell^2}{2t}\right).$$

Under this model, the rate of cooperation decreases gradually over time, in line with the experimental results of Figure 10(a) without delay. (The initial increase in cooperation rates in the fast treatments is likely due to subjects reacting relatively slowly to adjusting their action from the end of the previous block.) Of course, a richer model of passage times and cutoff strategies can be used to fit the data precisely by fitting an appropriate inverse Gaussian distribution to average cooperation rates and allowing for variation in thresholds  $\ell$ .

### 3.6 Discussion

Our experiments provide the first systematic treatment of delay, communication, and reaction lags in a repeated game with frequent actions and imperfect monitoring.

The results draw attention to several discrepancies with standard theory. The first finding is that delay unambiguously hurt subjects. This is at odds with the [Abreu et al. \(1991\)](#) argument that information delay can substantially help players to reduce the likelihood of inefficient punishments. The argument is robust—it holds in games with frequent and infrequent actions, with unbounded patience, even in approximate equilibrium, which makes the empirical finding all the more puzzling. One explanation of it is that delay makes it more difficult for subjects to learn the kind of opponent they are facing. Intuitively, if subjects face too much uncertainty over their opponents’ planned behavior, then it may be difficult to justify cooperating with them. The absence of delay may help players to signal their intended behavior more effectively.

Our second finding is that communication unambiguously helps subjects. Although it is well known in the experimental literature that communication generally improves welfare in a wide variety of strategic contexts, there is no strong theoretical justification for it adding (or subtracting) value in our experiment, unless—again—it helps to reduce subjects’ uncertainty over their opponents’ intentions. This suggests several sources of possible gains from reducing uncertainty over opponents’ strategies. First, it may motivate a subject to cooperate more if she is more confident that her opponent is likely to cooperate in return. Second, it may be easier to give incentives for cooperation to opponents if they can be made better aware of the consequences of their defection. In treatments with communication, where subjects announce their contingent plans at the beginning of every block, both of these channels should be facilitated—at least somewhat—and we find this to be the case. Subjects use messages to inform their opponents of future strategies, and reported and actual behaviors are significantly aligned.

We also find little evidence that reaction lags affect behavior in our experiment. This is a notable observation in light of the experiments reported in [Friedman and Oprea \(2012\)](#), where cooperation rates increase monotonically as players are given the opportunity to respond more quickly. [Friedman and Oprea \(2012\)](#) argued that response lags prevent players from quickly punishing deviations and that the gains to shorter lags are monotonic. Crucially, this argument only makes sense with *perfect* monitoring. In an environment like ours where information is noisy, reaction lags may allow players to gain a better idea of whether their partners are being cooperative, thereby supporting more efficient equilibria. We find, however, that cooperation does not seem to be affected in practice by how much time a player has to respond to a



signal.

Our instrumental variable-based estimation shows that subjects behavior is driven by observations of the public signal. Using an opponent's strategic behavior as an instrument is justified by the assumption that subjects play mixed strategies. As a result, conditional on the public signal, random changes in a subject's behavior must be mutually independent. Since the coefficient on the public signal is estimated with an instrument, it is unbiased regardless of other omitted variables. If there are no omitted variables then the regression equation is consistent with a first-order approximation of public strategies. Of course, if we assume that subjects' behavior is consistent with equilibrium then players must be playing public equilibria. From a practical point of view, public equilibrium makes testable restrictions on feasible outcomes, such as the welfare bound of [Proposition 1](#). However, our treatment with communication implies that the public equilibrium bound is violated. This could be for several reasons. First, it could be that subjects simply do not play equilibrium strategies. However, if the equilibrium assumption is dropped then it is not clear what structural predictions can be made. Secondly, it could be that their bounded rationality means that they are incapable of playing public equilibria, as they cannot react immediately to an individual bad news event. This may improve welfare, as illustrated in [Section 3.4.4](#). However, it is not clear why players would exploit this bounded rationality in treatments with communication in order to exceed the public equilibrium bound but not in those without. Thirdly, subjects may be playing private—not public—strategies. From a technical point of view, public equilibrium is an assumption that often puts severe restrictions on behavior. That is, it precludes players from certain behavior that may be intuitively justified in some contexts. [Rahman \(2012, 2013a\)](#) explores this issue at some length and argues that public equilibria preclude secret monitoring and infrequent coordination amongst players. Both of these behaviors have the potential to improve their welfare significantly, so much so that the impossibility results of [Sannikov and Skrzypacz \(2007, 2010\)](#) can be completely overturned. Our interest in future work is to explore experimentally how infrequent coordination can help sustain cooperation in the laboratory.

## References

- Abreu, D. and F. Gul (2000): “Bargaining and reputation,” *Econometrica*, 68, 85–117. 28
- Abreu, D., P. Milgrom, and D. G. Pearce (1991): “Information and Timing in Repeated Partnerships,” *Econometrica*, 59, 1713–33. 49, 52, 55, 59, 60, 62, 63, 77
- Abreu, D., D. Pearce, and E. Stacchetti (1986): “Optimal cartel equilibria with imperfect monitoring,” *Journal of Economic Theory*, 39, 251–269. 55
- Abreu, D., D. G. Pearce, and E. Stacchetti (1990): “Toward a Theory of Discounted Repeated Games with Imperfect Monitoring,” *Econometrica*, 58, 1041–1063. 55
- Alchian, A. and H. Demsetz (1972): “Production, information costs, and economic organization,” *The American Economic Review*, 62, 777–795. 26
- Almlund, M., A. L. Duckworth, J. J. Heckman, and T. Kautz (2011): “Personality psychology and economics,” in *Handbook of the economics of education*, ed. by E. A. Hanushek, S. Machin, and L. Woessmann, Amsterdam: Elsevier, 1–181. 26
- Ambrus, A. and B. Greiner (2012): “Imperfect Public Monitoring with Costly Punishment: An Experimental Study,” *The American Economic Review*, 102, 3317–3332. 53
- Anderson, J., S. Burks, C. D. C, and A. Rustichini (2011): “Toward the integration of personality theory and decision theory in the explanation of economic behavior,” Unpublished manuscript. 27, 29, 30, 96
- Angrist, J. D. and J.-S. Pischke (2008): *Mostly harmless econometrics: An empiricist’s companion*, Princeton university press. 41
- Aoyagi, M. and G. Fréchette (2009): “Collusion as public monitoring becomes noisy: Experimental evidence,” *Journal of Economic theory*, 144, 1135–1165. 53
- Armantier, O. and N. Treich (2013): “Eliciting beliefs: Proper scoring rules, incentives, stakes and hedging,” *European Economic Review*, 62, 17–40. 5
- Ashton, M. C., K. Lee, L. R. Goldberg, and R. E. de Vries (2009): “Higher order factors of personality: do they exist?” *Personality and Social Psychology Review*, 13, 79–91. 96
- Bagwell, K. and G. Ramey (1996): “Capacity, Entry, and Forward Induction,” *The Rand Journal of Economics*, 660–680. 1, 4
- Balkenborg, D. (1994): *An Experiment on Foreward Versus Backward Induction*, Sonderforschungsbereich 303. 1, 4

- Battigalli, P. and A. Friedenberg (2012): “Forward induction reasoning revisited,” *Theoretical Economics*, 7, 57–98. 1, 4
- Battigalli, P. and M. Siniscalchi (2002): “Strong Belief and Forward Induction Reasoning,” *Journal of Economic Theory*, 106, 356–391. iii, 4
- Becker, A., T. Deckers, T. Dohmen, A. Falk, and F. Kosse (2012): “The relationship between economic preferences and psychological personality measures,” *Annual Review of Economics*, 4, 45378. 27, 30, 47
- Ben-Ner, A. and A. Kramer (2011): “Personality and altruism in the dictator game: Relationship to giving to kin, collaborators, competitors, and neutrals,” *Personality and Individual Differences*, 51, 216–221. 29, 30
- Ben-Ner, A., A. Kramer, and O. Levy (2008): “Economic and hypothetical dictator game experiments: Incentive effects at the individual level,” *The Journal of Socio-Economics*, 37, 1775 – 1784. 30
- Ben-Ner, A., L. Putterman, and T. Ren (2007): “Lavish returns on cheap talk: non-binding communication in a trust experiment,” Working paper. 54
- Ben-Porath, E. and E. Dekel (1992): “Signaling future actions and the potential for sacrifice,” *Journal of Economic Theory*, 57, 36–51. 1, 4
- Biddle, J. and D. Hamermesh (1998): “Beauty, productivity, and discrimination: Lawyers’ Looks and Lucre,” *Journal of Labor Economics*, 16, 172–201. 99
- Bigoni, M., M. Casari, A. Skrzypacz, and G. Spagnolo (2011): “Time horizon and cooperation in continuous time,” Working paper. 53
- Bigoni, M., J. J. Potters, and G. Spagnolo (2012): “Flexibility and Collusion with Imperfect Monitoring,” Working paper. 53
- Binmore, K., A. Rubinstein, and A. Wolinsky (1986): “The Nash bargaining solution in economic modelling,” *The RAND Journal of Economics*, 176–188. 28
- Blanco, M., D. Engelmann, A. Koch, and H.-T. Normann (2010): “Belief Elicitation in Experiments: Is There a Hedging Problem?” *Experimental Economics*, 13, 412–438. 23
- Blume, A. and U. Gneezy (2010): “Cognitive forward induction and coordination without common knowledge: An experimental study,” *Games and Economic Behavior*, 68, 488–511. 4

- Bohnet, I. and B. Frey (1999): “Social distance and other-regarding behavior in dictator games: comment,” *American Economic Review*, 89, 335–339. 32
- Borghans, L., A. Duckworth, J. Heckman, and B. Ter Weel (2008): “The economics and psychology of personality traits,” *Journal of Human Resources*, 43, 972–1069. 27, 30
- Bouchard, ThomasJ., J. and J. Loehlin (2001): “Genes, evolution, and personality,” *Behavior Genetics*, 31, 243–273. 99
- Brandstätter, H. and M. Königstein (2001): “Personality influences on ultimatum bargaining decisions,” *European Journal of Personality*, 15, S53–S70. 29
- Brandts, J., A. Cabrales, and G. Charness (2007): “Forward induction and entry deterrence: an experiment,” *Economic Theory*, 33, 183–209. 1, 4
- Brandts, J. and C. A. Holt (1995): “Limitations of dominance and forward induction: Experimental evidence,” *Economics Letters*, 49, 391–395. 1, 4
- Burke, P. (2003): “Interaction is small groups,” in *Handbook of Social Psychology*, ed. by J. DeLamater, New York: Kluwer-Plenum, 363–388. 31
- Cachon, G. P. and C. F. Camerer (1996): “Loss-avoidance and forward induction in experimental coordination games,” *The Quarterly Journal of Economics*, 111, 165–194. 4
- Camerer, C. and R. Hogarth (1999): “The effects of financial incentives in experiments: A review and capital-labor-production framework,” *Journal of risk and uncertainty*, 19, 7–42. 30
- Cason, T. N. and F. U. Khan (1999): “A laboratory study of voluntary public goods provision with imperfect monitoring and communication,” *Journal of Development Economics*, 58, 533–552. 53
- Charness, G. (2000): “Self-serving cheap talk: A test of Aumann’s conjecture,” *Games and Economic Behavior*, 33, 177–194. 53, 54
- Charness, G. and M. Dufwenberg (2006): “Promises and Partnership,” *Econometrica*, 74, 1579–1601. 6, 11
- Charness, G., R. Oprea, and D. Friedman (2012): “Continuous Time and Communication in a Public-goods Experiment,” Working paper. 54
- Cho, I.-K. and D. M. Kreps (1987): “Signaling games and stable equilibria,” *The Quarterly Journal of Economics*, 102, 179–221. 4

- Chung, K.-S. (2011): “Forward Induction and Changing Tastes,” Mimeo. 4
- Clark, J. (1902): *The distribution of wealth: a theory of wages, interest and profit*, The Macmillan Company. 26
- Compte, O. (1998): “Communication in Repeated Games with Imperfect Private Monitoring,” *Econometrica*, 66, 597–626. 49, 50, 55, 59
- Cooper, R., D. DeJong, R. Forsythe, and T. Ross (1993): “Forward Induction in the Battle-of-the-Sexes Games,” *American Economic Review*, 83, 1303–1316. 1, 4, 5, 14, 15
- Costa, P. T. and R. R. McCrae (1992): *Revised NEO personality inventory (NEO-PI-R) and NEO five-factor inventory (NEO-FFI) manual*, Psychological Assessment Resources. iv, 27
- Costa-Gomes, M. A. and G. Weizsäcker (2008): “Stated beliefs and play in normal-form games,” *The Review of Economic Studies*, 75, 729–762. 5
- Crawford, V. (1998): “A survey of experiments on communication via cheap talk,” *Journal of Economic theory*, 78, 286–298. 32, 50
- Crawford, V. P. and J. Sobel (1982): “Strategic information transmission,” *Econometrica: Journal of the Econometric Society*, 1431–1451. 50
- Croson, R. T. (2000): “Thinking like a game theorist: factors affecting the frequency of equilibrium play,” *Journal of economic behavior & organization*, 41, 299–314. 5
- Dal Bó, P. (2005): “Cooperation under the shadow of the future: experimental evidence from infinitely repeated games,” *The American Economic Review*, 95, 1591–1604. 53
- Dana, J., R. Weber, and J. Kuang (2007): “Exploiting moral wiggle room: experiments demonstrating an illusory preference for fairness,” *Economic Theory*, 33, 67–80. 32
- Danz, D. N., D. Fehr, and D. Kübler (2012): “Information and beliefs in a repeated normal-form game,” *Experimental Economics*, 15, 622–640. 5
- Dawes, R. M., J. McTavish, and H. Shaklee (1977): “Behavior, communication, and assumptions about other people’s behavior in a commons dilemma situation.” *Journal of personality and social psychology*, 35, 1. 53
- De Sinopoli, F. (2004): “A note on forward induction in a model of representative democracy,” *Games and Economic Behavior*, 46, 41–54. 4

- DeYoung, C., J. B. Hirsh, M. S. Shane, X. Papadermetris, and J. Gray (2010): “Testing predictions from personality neuroscience: brain structure and the big five,” *Psychological Science*, 21, 820–828. 27
- DeYoung, C. G., L. C. Quilty, and J. B. Peterson (2007): “Between facets and domains: 10 aspects of the big five,” *Journal of Personality and Social Psychology*, 93. 33, 39
- Edgeworth, F. (1881): *Mathematical psychics: An essay on the application of mathematics to the moral sciences*, 10, CK Paul. 26
- Ely, J. C., J. Hörner, and W. Olszewski (2005): “Belief-Free Equilibria in Repeated Games,” *Econometrica*, 73, 377–415. 49, 55, 59
- Feuerstein, S. (2005): “Collusion in industrial economics—a survey,” *Journal of Industry, Competition and Trade*, 5, 163–198. 52
- Fischbacher, U. (2007a): “z-Tree: Zurich Toolbox for Ready-made Economic Experiments,” *Experimental Economics*, 10, 171–178. 6, 33
- (2007b): “z-Tree: Zurich toolbox for ready-made economic experiments,” *Experimental economics*, 10, 171–178. 55
- Friedman, D. and R. Oprea (2012): “A continuous dilemma,” *The American Economic Review*, 102, 337–363. 51, 52, 53, 54, 77
- Fudenberg, D. and D. Levine (2007): “Continuous time limits of repeated games with imperfect public monitoring,” *Review of Economic Dynamics*, 10, 173–192. 55
- (2009): “Repeated Games with Frequent Signals,” *The Quarterly Journal of Economics*, 124, 233–265. 55
- Fudenberg, D., D. Levine, and E. Maskin (1994): “The Folk Theorem with Imperfect Public Information,” *Econometrica*, 62, 997–1039. 49, 55, 113
- Geanakoplos, J., D. Pearce, and E. Stacchetti (1989): “Psychological games and sequential rationality,” *Games and Economic Behavior*, 1, 60–79. 28
- Gill, D. and V. Prowse (2012): “A structural analysis of disappointment aversion in a real effort competition,” *American Economic Review*, 102, 469–503. 33
- Glazer, J. and A. Weiss (1990): “Pricing and coordination: Strategically stable equilibria,” *Games and Economic Behavior*, 2, 118–128. 4
- Goldberg, L. (1993): “The structure of phenotypic personality traits,” *American Psychologist*, 48, 26–34. 27

- Govindan, S. and R. Wilson (2009): “On forward induction,” *Econometrica*, 77, 1–28. 1, 4
- Gul, F. and D. G. Pearce (1996): “Forward induction and public randomization,” *journal of economic theory*, 70, 43–64. 4
- Hauk, E. and S. Hurkens (2002): “On forward induction and evolutionary and strategic stability,” *Journal of Economic Theory*, 106, 66–90. 4
- Heckman, J., J. Stixrud, and S. Urzua (2006): “The effects of cognitive and noncognitive abilities on labor market outcomes and social behavior,” Tech. rep., National Bureau of Economic Research. 29, 30
- Holmstrom, B. (1982): “Moral hazard in teams,” *The Bell Journal of Economics*, 324–340. 26
- Hörner, J. and W. Olszewski (2006): “The Folk Theorem for Games with Private Almost-Perfect Monitoring,” *Econometrica*, 74, 1499–1544. 49, 59
- Houser, D. and E. Xiao (2011): “Classification of natural language messages using a coordination game,” *Experimental Economics*, 14, 1–14. 6
- Huck, S. and W. Müller (2005): “Burning Money and (Pseudo) First-mover Advantages: An Experimental Study on Forward Induction,” *Games and Economic Behavior*, 51, 109–127. 1, 4
- Hyndman, K., E. Y. Ozbay, A. Schotter, and W. Z. Ehrblatt (2012): “Convergence: an experimental study of teaching and learning in repeated games,” *Journal of the European Economic Association*, 10, 573–604. 5
- Judge, T. A., B. A. Livingston, and C. Hurst (2012): “Do nice guys—and gals—really finish last? The joint effects of sex and agreeableness on income,” *Journal of Personality and Social Psychology*, 102, 390–407. 30
- Kahneman, D. (2011): “Thinking, Fast and Slow,” *New York: Farrar, Straus and Giroux*. 7
- Kandori, M. (2003): “Randomization, Communication, and Efficiency in Repeated Games with Imperfect Public Monitoring,” *Econometrica*, 71, 345–353. 50
- Kandori, M. and H. Matsushima (1998): “Private Observation, Communication, and Collusion,” *Econometrica*, 66, 627–652. 50, 55, 59, 63
- Kandori, M. and I. Obara (2006): “Efficiency in Repeated Games Revisited: The Role of Private Strategies,” *Econometrica*, 74, 499–519. 49

- Knight, F. (1921): “Risk, uncertainty and profit,” *New York: Hart, Schaffner and Marx*. 26, 30, 31
- Kohlberg, E. and J.-F. Mertens (1986): “On the Strategic Stability of Equilibria,” *Econometrica*, 54, 1003–1037. iii, 1, 2, 4
- Kreps, D. M., P. Milgrom, J. Roberts, and R. Wilson (1982): “Rational cooperation in the finitely repeated prisoners’ dilemma,” *Journal of Economic theory*, 27, 245–252. 54
- Kyl-Heiku, L. and D. Buss (1996): “Tactics as units of analysis in personality psychology: An illustration using tactics of hierarchy negotiation,” *Personality and Individual Differences*, 21, 497–517. 47
- Lehrer, E. (1989): “Lower equilibrium payoffs in two-player repeated games with non-observable actions,” *International Journal of Game Theory*, 18, 57–89. 49, 55, 60
- Levin, J. (2003): “Relational Incentive Contracts,” *American Economic Review*, 93, 835–847. 26, 52
- Man, P. T. (2012): “Forward induction equilibrium,” *Games and Economic Behavior*, 75, 265–276. 1
- Mobius, M. and T. Rosenblat (2006): “Why beauty matters,” *American Economic Review*, 96, 222–235. 29, 99
- Monson, T., J. Hesley, and L. Chernick (1982): “Specifying when personality traits can and cannot predict behavior: An alternative to abandoning the attempt to predict single-act criteria.” *Journal of Personality and Social Psychology*, 43, 385. 31
- Morris, M., R. Larrick, and S. Su (1999): “Misperceiving negotiation counterparts: When situationally determined bargaining behaviors are attributed to personality traits.” *Journal of Personality and Social Psychology*, 77, 52. 29, 30, 31, 38
- Mueller, G. and E. Plug (2006): “Estimating the effect of personality on male and female earnings,” *Industrial and Labor Relations Review*, 60, 3–22. 30
- Müller, C. (2012): “Robust virtual implementation under common strong belief in rationality,” Mimeo. 1, 4
- Muller, R. A. and A. Sadanand (2003): “Order of play, forward induction, and presentation effects in two-person games,” *Experimental Economics*, 6, 5–25. 4
- Murnighan, J. K. and A. E. Roth (1983): “Expecting Continued Play in Prisoner’s Dilemma Games A Test of Several Models,” *Journal of Conflict Resolution*, 27, 279–300. 56



- Nadler, J. and D. Shestowsky (2006): “Negotiation, information technology, and the problem of the faceless other,” in *Negotiation Theory and Research*, ed. by L. L. Thompson, New York: Psychology Press. 32
- Nash Jr, J. (1950): “The bargaining problem,” *Econometrica: Journal of the Econometric Society*, 155–162. 28
- Ng, T. W. H., L. T. Eby, K. L. Soresen, and D. C. Feldman (2005): “Predictors of objective and subjective career success: A meta-analysis,” *Personnel Psychology*, 58, 367–408. 30
- Nyarko, Y. and A. Schotter (2002): “An experimental study of belief learning using elicited beliefs,” *Econometrica*, 70, 971–1005. 5
- Nyhus, E. K. and E. Pons (2005): “The effects of personality on earnings,” *Journal of Economic Psychology*, 26, 363 – 384. 30
- Obara, I. (2008): “Folk Theorem with Communication,” *Journal of Economic Theory*, forthcoming, mimeo. 55
- (2009): “Folk theorem with communication,” *Journal of Economic Theory*, 144, 120–134. 50
- Oprea, R., K. Henwood, and D. Friedman (2011): “Separating the Hawks from the Doves: Evidence from continuous time laboratory games,” *Journal of Economic Theory*, 146, 2206–2225. 53
- Papke, L. and J. Woolridge (2008): “Panel data methods for fractional response variables with an application to test pass rates,” *Journal of Econometrics*, 145, 121–133. 45
- Ponssard, J.-P. (1991): “Forward induction and sunk costs give average cost pricing,” *Games and Economic Behavior*, 3, 221–236. 1, 4
- Porter, R. H. (2005): “Collusion in industrial economics: a comment,” *Journal of Industry, Competition and Trade*, 5, 231–234. 52
- Pound, N., I. S. Penton-Voak, and W. M. Brown (2007): “Facial symmetry is positively associated with self-reported extraversion,” *Personality and Individual Differences*, 43, 1572 – 1582. 99
- Prendergast, C. (1999): “The provision of incentives in firms,” *Journal of economic literature*, 37, 7–63. 26
- Radner, R. (1986): “Can bounded rationality resolve the prisoners dilemma?” *Essays in honor of Gerard Debreu*, 387–399. 51, 54

- Radner, R., R. Myerson, and E. Maskin (1986): “An Example of a Repeated Partnership Game with Discounting and with Uniformly Inefficient Equilibria,” *Review of Economic Studies*, 53, 59–69. 55
- Rahman, D. (2012): “Mediating Collusion with Flexible Production,” Mimeo. 78
- (2013a): “Frequent Actions with Infrequent Coordination,” Working paper. 50, 78, 113
- (2013b): “Information Delay in Games with Frequent Actions,” Working paper. 55, 63
- Rick, S. and R. A. Weber (2010): “Meaningful Learning and Transfer of Learning in Games Played Repeatedly Without Feedback,” *Games and Economic Behavior*, 68, 716–730. 7
- Rode, J. C., M. L. Arthaud-Day, C. H. Mooney, J. P. Near, and T. T. Baldwin (2008): “Ability and Personality Predictors of Salary, Perceived Job Success, and Perceived Career Success in the Initial Career Stage,” *International Journal of Selection and Assessment*, 16, 292–299. 30
- Rubinstein, A. (1982): “Perfect equilibrium in a bargaining model,” *Econometrica: Journal of the Econometric Society*, 97–109. 28
- Rushton, J. P. and P. Irwing (2008): “A General Factor of Personality (GFP) from two meta-analyses of the Big Five: and,” *Personality and Individual Differences*, 45, 679 – 683. 96
- Samuelson, L. (2005): “Economic Theory and Experimental Economics,” *Journal of Economic Literature*, 43, 65–107. 1
- Sannikov, Y. (2007): “Games with imperfectly observable actions in continuous time,” *Econometrica*, 75, 1285–1329. 55
- Sannikov, Y. and A. Skrzypacz (2007): “Impossibility of collusion under imperfect monitoring with flexible production,” *The American Economic Review*, 97, 1794–1823. 49, 55, 78
- (2010): “The role of information in repeated games with frequent actions,” *Econometrica*, 78, 847–882. 55, 78
- Schotter, A. and B. Sopher (2007): “Advice and behavior in intergenerational ultimatum games: An experimental approach,” *Games and Economic Behavior*, 58, 365–393. 6

- Schotter, A. and I. Trevino (2014): “Belief Elicitation in the Laboratory,” *Annual Review of Economics*. 5
- Seibert, S. and M. Kraimer (2001): “The five-factor model of personality and career success,” *Journal of vocational behavior*, 58, 1–21. 29
- Selten, R., A. Sadrieh, and K. Abbink (1999): “Money Does Not Induce Risk Neutral Behavior, but Binary Lotteries Do Even Worse,” *Theory and Decision*, 46, 213–252. 14
- Shahriar, Q. (2013): “An Experimental Test of the Robustness and the Power of Forward Induction,” *Managerial and Decision Economics*. 4
- Simon, L. K. and M. B. Stinchcombe (1989): “Extensive form games in continuous time: Pure strategies,” *Econometrica: Journal of the Econometric Society*, 1171–1214. 51, 54
- Sonnemans, J. and T. T. Offerman (2001): “Is the quadratic scoring rule behaviorally incentive compatible?” Mimeo. 5
- Sugaya, T. (2010): “Folk Theorem in Repeated Games with Private Monitoring,” Tech. rep., mimeo. 55, 59
- Teglas, E., E. Vul, V. Girotto, M. Gonzalez, J. B. Tenenbaum, and L. L. Bonatti (2011): “Pure Reasoning in 12-Month-Old Infants as Probabilistic Inference,” *Science*, 332, 1054–1059. 25
- Thompson, L. (1990): “Negotiation behavior and outcomes: Empirical evidence and theoretical issues.” *Psychological bulletin*, 108, 515. 31
- Thompson, L., J. Wang, and B. Gunia (2010): “Negotiation,” *Annual review of psychology*, 61, 491–515. 29
- Ulen, T. S. (1980): “Cartels and Regulation: Late Nineteenth-Century Railroad Collusion and the Creation of the Interstate Commerce Commission,” *The Journal of Economic History*, 40, pp. 179–181. 52
- Van Damme, E. (1989): “Stable equilibria and forward induction,” *journal of Economic Theory*, 48, 476–496. 4
- Van Huyck, J. B., R. C. Battalio, and R. O. Beil (1993): “Asset markets as an equilibrium selection mechanism: Coordination failure, game form auctions, and tacit communication,” *Games and Economic Behavior*, 5, 485–504. 4
- Willis, J. and A. Todorov (2006): “First impressions: making up your mind after a 100-ms exposure to a face,” *Psychological Science*, 17, 592–598. 33

Xiao, E. and D. Houser (2005): "Emotion expression in human punishment behavior,"  
*Proceedings of the National Academy of Sciences of the United States of America*, 102,  
7398–7401. 2, 6

## 4 Appendix

### 4.1 Questions for Second Movers

While first movers were answering F1-F3, second movers answered S1-S4.

S1. What will the first mover do?

Possible answers: *In* for sure, 90/10 *In/Out*, 80/20 *In/Out*, ..., *Out* for sure.

S2. If the first mover goes in, what will he/she think about the move you will make in the second stage?

Possible answers: *Left* for sure, 90/10 *Left/Right*, 80/20 *Left/Right*, ..., *Right* for sure.

S3. If the first mover goes in, what will he/she do next?

Possible answers: *Top* for sure, 90/10 *Top/Bottom*, 80/20 *Top/Bottom*, ..., *Bottom* for sure.

After the choice of the first mover was observed, all second movers answered S4.

S4. The first mover thought that if the second mover is given a chance to move, the second mover will pick:

Possible answers: *Left* for sure, 90/10 *Left/Right*, 80/20 *Left/Right*, ..., *Right* for sure.

The remaining questions of second movers depended on whether their partner chose *Out* or *In*. The choice of *Out* prompted S5-S7.

S5. Why did the first mover think this? (This question followed S4.)

S6. If you were given a chance to move, you would choose:

Possible answers: *Left* for sure, 90/10 *Left/Right*, 80/20 *Left/Right*, ..., *Right* for sure.

S7. Why would you make this move?

If the first mover chose *In*, the second mover answered S8-S10.

S8. For his/her next move, the first mover will choose:

Possible answers: *Top* for sure, 90/10 *Top/Bottom*, 80/20 *Top/Bottom*, ..., *Bottom* for sure.

S9. Why will the first mover make this choice?

S10. Why did you make this move? (This question followed the second mover's choice between *Left* and *Right*.)

## 4.2 Instructions

Welcome to this decision making experiment.

### INSTRUCTIONS

**You will be quizzed on these instructions, and you will not be able to continue the experiment until you complete the quiz. Please read the instructions carefully, and raise your hand if you have any questions.**

In this experiment, you will play a game a number of times. The whole experiment (including reading the instructions) should take no longer than 1 hour.

In this game, there is a first mover and a second mover. Depending on what the first mover does, **one round of this game can have 1 or 2 stages.**

#### ONE ROUND

##### **STAGE 1 (happens always)**

In this stage, the first mover will decide whether to choose Out or In.

- If the first mover chooses Out, both players will get \$7. The second mover will be informed that the first mover went Out, and the round will end.
- If the first mover chooses In, the game will proceed to the next stage.

##### **STAGE 2 (happens only in the first mover chooses In)**

In this stage, the second mover will be informed that the first mover went In, and both players will make additional choices. The first mover will choose between Top and Bottom and the second mover will choose between Left and Right.

- If the first mover chooses Top, and the second mover chooses Left, the first mover will get \$10 and the second mover will get \$5.
- If the first mover chooses Top, and the second mover chooses Right, both players will get \$0.
- If the first mover chooses Bottom, and the second mover chooses Left, both players will get \$0.
- If the first mover chooses Bottom, and the second mover chooses Right, the first mover will get \$5 and the second mover gets \$10.
- After both players make their choices, the round will end.

**IMPORTANT: Neither player will be informed of what the other person did in Stage 2. Therefore, the first mover will not find out what the second mover did in Stage 2, and the second mover will not find out what the first mover did in Stage 2.**

- After the round ends, a new round starts.

Go on to the next page for more instructions

## YOUR ROLE IN THE GAME

Half of the participants will start out as first movers, and half will start out as second movers. Your role will switch throughout the experiment. Therefore, if you start out as the first mover, your roles will be: first mover (round 1), second mover (round 2), first mover (round 3), second mover (round 4), etc.

You will never know the identity of the person you are playing with.

## HOW YOU ARE MATCHED WITH OTHER PEOPLE

In each round, you will be randomly matched with another player. Because the matching is random, it is very unlikely that you will be matched with the same person twice in a row.

## OTHER TASKS YOU HAVE TO COMPLETE

We will ask you questions about your thoughts regarding the game. Some of the answers are “multiple choice,” and some are “short answer” (you have to type).

You will type the “short answers” in the blue box. **PRESS THE ENTER KEY WHEN YOU ARE DONE TYPING.** Otherwise, your answer will not save. If you want to type more after you hit ENTER, type in the blue box again, and press the ENTER key again.

Go on to the next page for information on how you get paid



## HOW YOU GET PAID

1. **You will be paid for your choices in the games. You will get this money immediately at the end of the experiment.**

We will randomly select 2 rounds of the game you played. If you got X dollars in the first round and Y dollars in the second round, we will pay you X+Y dollars for the games.

2. **You will later be paid for your answers to the questions we ask.**

You will see questions, in which you will be

- guessing the other player's behavior,
- guessing the other player's thoughts,
- explaining the other person's behavior or thoughts.

When we analyze your answers to these questions, we will score them against the response of the other player. **It is always in your best interests to answer them as truthfully as possible.** The closer your answer is to the truth (which we get from the other player), the more money you make. Within 3 weeks from today, we will mail you a check for your guesses and explanations.

We will not pay you additional money for answers to questions about **your own** thoughts and behavior, but obtaining your considered opinion is important for our study, so please be as detailed as possible in your answers.

The following page of the instructions describes the details of how payments for your answers are calculated. You will not be quizzed on this information, but skim it if you need to be convinced that answering these questions truthfully is in your best interests.

**DETAILS: Payments for short answers**

For each question in which you explain what the other player is doing (or thinking), this player will have a corresponding question in which they have to explain their own behavior (or thoughts). We will evaluate your answer by placing it into one of several categories, and do the same for the answer the other player provides. If the categories match, we will pay you \$1. For example, if your explanation of why the other player will behave in a particular way matches that player's explanation of their behavior, you get \$1. If your explanations do not match, we will pay you \$0. If one of you does not provide a short answer, we will pay you \$0.

**DETAILS: Payments for multiple choice questions**

When you are making a guess about the other player's **behavior**, the amount you are paid will be calculated as in the example below.

**Example:** The other person chooses between Move A and Move B, and you guess what move they will make. Your possible guesses – the odds of the other player making one choice or the other – are in the leftmost column of the table below. Your payments, which depend on the other player's choice, are in the other two columns. **The closer your guess is to their choice, the more money you make.**

	Other person chooses Move A	Other person chooses Move B
You guess "Move A for sure"	100 cents	0
You guess "90/10 Move A/Move B"	99 cents	19 cents
You guess "80/20 Move A/Move B"	96 cents	36 cents
You guess "70/30 Move A/Move B"	91 cents	51 cents
You guess "60/40 Move A/Move B"	84 cents	64 cents
You guess "50/50 Move A/Move B"	75 cents	75 cents
You guess "40/60 Move A/Move B"	64 cents	84 cents
You guess "30/70 Move A/Move B"	51 cents	91 cents
You guess "20/80 Move A/Move B"	36 cents	96 cents
You guess "10/90 Move A/Move B"	19 cents	99 cents
You guess "Move B for sure"	0 cents	100 cents

When you are making a guess about the other player's **thoughts**, the answers you and the other player provide will be probabilities, and your payment will be calculated as  $1 - (p - q)^2$ , where  $p$  is the probability you reported, and  $q$  is the probability reported by the other player. **The closer your report is to the report of the other player, the more money you make.**

We will pay you for your multiple choice answers whenever possible. Some of these questions you will answer will be of the form "if *something*, then [your multiple choice guess]." Whenever this "something" is not realized, we will not be able to score your answers. But we might still use them when we are analyzing your written explanations! **It is always in your best interests to answer multiple choice questions truthfully.**

If you guess about the same item more than once in the course of one round, and both answers are scored, then we will pay you for your best guess.

### 4.3 Quiz questions for participants

- You are the first mover, and you go Out. How much does the second mover get? (a) \$0; (b) \$7; (c) \$5; (d) \$10.
- You are the second mover. The first mover went In and chose Bottom for their next move. How much do you get? (a) \$0; (b) \$7; (c) If I go Left, I get \$0. If I go Right, I get \$10; (d) \$5.
- In this game, the first mover will be informed of: (a) (If he/she goes In) Whether the second mover chose Left or Right; (b) (If he/she goes In) How much money both players made; (c) The first mover will not be informed of anything that happens after he/she goes In; (d) None of the above is true.
- In this game, the second mover will be informed of: (a) (If the first mover went In) Whether the first mover chooses Top or Bottom next; (b) Whether the first mover went In or Out; (c) The second mover will not be informed of anything that happens in this game; (d) None of the above is true.

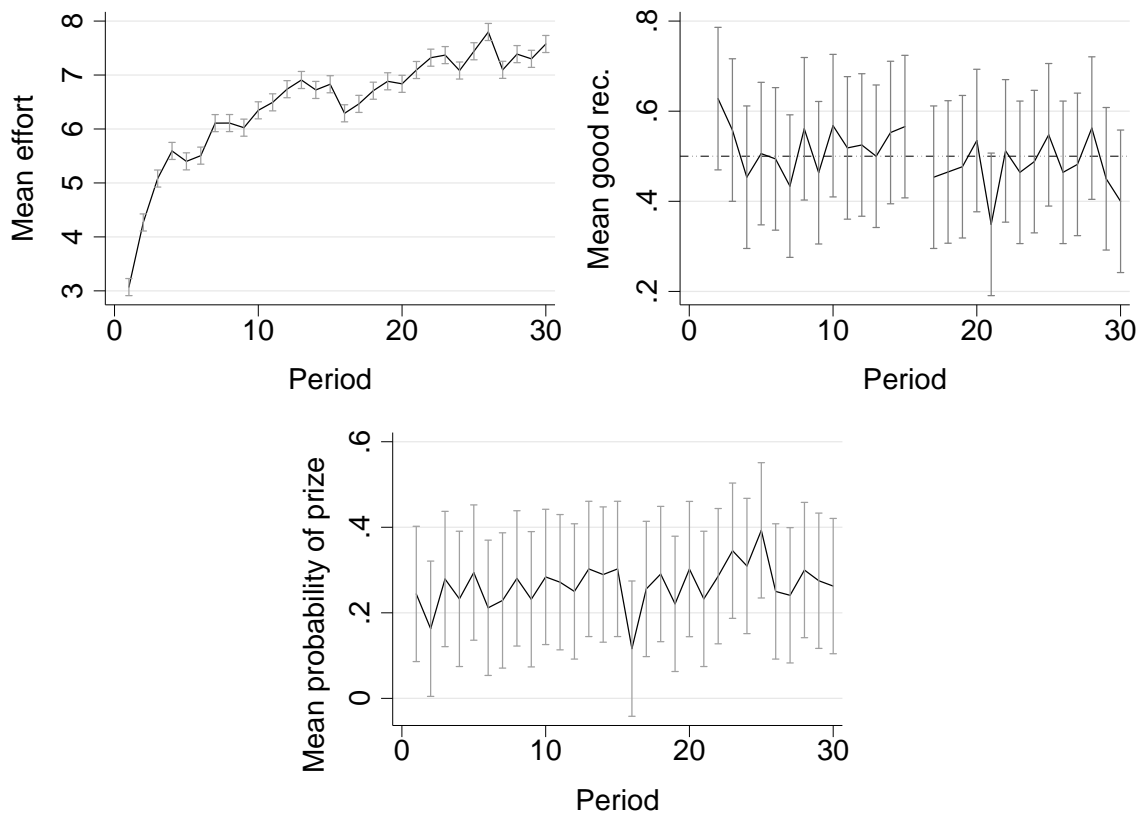
### 4.4 Summary of the Data

There was substantial variation in self-reported personality traits. Although we found significant correlations between all traits in the Big Five (Table 18), the signs were consistent with what has been reported in the literature (Anderson et al, 2011).<sup>45</sup> Cronbach’s alphas, which provide a reliability measure for how good a set of questions is at capturing a particular trait, are high in our data, ranging from 0.85 for Neuroticism to 0.9 for Extraversion.

Figure 12a shows how the amount of correctly adjusted sliders evolved over time. The workers moved close to seven sliders on average, but there was substantial learning, especially in the early portion of the experiment. This was confirmed by an OLS regression of the number of sliders moved in the course of the match against a half dummy ( $P < 0.001$ ; see Figure 12d). Figure 12b shows the dynamics of managers’ recommendations, which on average were not better than chance. An OLS regression of the number of good recommendations against a half dummy showed that the recommendations were significantly worse in the second half of the experiment ( $P < 0.001$ ). The average probability of discovering a prize over time is plotted in Figure 12c. The teams did substantially worse in the first period after teams were re-matched. Likely, this was the case because the prize was behind

---

<sup>45</sup>It is well-known that personality traits are correlated, although the reasons for such correlations are unclear. Evidence of “meta-traits” has been documented (Rushton and Irwing, 2008), but the topic remains controversial (Ashton et al, 2009).



	No. of sliders moved	No. of good recommendations	No. of prizes discovered
Half	17.78**** (1.815)	-0.797** (0.266)	0.131 (0.219)
Constant	72.52**** (4.916)	8.245**** (0.440)	4.001**** (0.428)
Observations	156	156	156

Standard errors in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ , \*\*\*\*  $p < 0.001$

Figure 12: Sliders moved (a), good recommendations made (b), and prizes discovered (c) over time. Regressions of these variables against a half dummy show that subjects moved more sliders in the second half of the experiment, that managers made worse recommendations in the second half, and that the amount of prizes discovered did not significantly differ across halves (d).

	N	A	C	E	O
N	1.0000				
A	-0.1671	1.0000			
C	-0.3362	0.2797	1.0000		
E	-0.3063	0.2582	0.4174	1.0000	
O	-0.2970	0.4037	0.2386	0.4280	1.0000

Table 18: Correlations between self-reported personality traits.

a bottom slider in this period: Starting out by moving top sliders was more natural for the workers. There was, however, no significant half effect on the number of prizes discovered ( $P = 0.5654$ ). On average, the probability of discovering the prize in any period of the interaction was 0.26. Thus, over the course of 15 periods, close to four prizes were accumulated, leaving the manager with approximately \$16 to divide between herself and the worker.



Figure 13: Money that the manager allocated to the worker and herself, as a function of whether or not a prize was discovered. Period number is on the horizontal axis.

Contrary to the non-cooperative prediction, managers did not appropriate the entire accumulated earnings in the last period. Figure 13 plots the amounts of money that the manager allocated to herself and her worker in every period, as a function of whether or

not a prize was discovered. Both parties made more money when a prize was discovered than when it wasn't. In the latter case, the amounts allocated to the worker were positive but small, and the manager's earnings negative due to the 40 cent continuation fee. The figures suggest that amounts allocated per period did not change over time.

### *Beauty and gender*

There is a considerable literature on the relationship between personality and gender, in which it has been reported, e.g., in a review by [Bouchard and Loehlin \(2001\)](#), that women score higher in agreeableness (and neuroticism) than men. Although we don't know of any research linking agreeableness to beauty, there is evidence of a relationship between facial symmetry and self-reported extraversion ([Pound et al, 2007](#)). Thus, a link between beauty and personality cannot be ruled out. On the other hand, earnings gaps due to gender and beauty are well-established, raising the possibility of agreeableness being correlated with the error in our regressions.

To address these concerns, we photographed participants in eight of the ten sessions we ran at the end of the experiment.<sup>46</sup> This allowed us to control for their gender as well as beauty in this subsample. We hired 20 University of Minnesota undergraduates to rate the photos for a flat fee of \$10. The same subject pool was used, but students who participated in the original experiment were not allowed to sign up for the photo evaluation component. All raters were presented with the entire set of photographs (order randomized) and told to rate each photo on a scale of 1-5, borrowed from [Biddle and Hamermesh \(1998\)](#).<sup>47</sup> We followed the procedure described in [Mobius and Rosenblat \(2006\)](#) to create our beauty regressor: rater  $j$ 's centered beauty rating  $\tilde{r}_{ji}$  of subject  $i$  was obtained by subtracting the rater's average beauty rating  $\bar{r}_j$  from each raw rating  $r_{ji}$ , and the variable  $Beauty_i$  defined as the mean over  $j$  of  $\tilde{r}_{ji}$ .

## 4.5 Additional Robustness Checks

---

<sup>46</sup>We were not able to obtain the photographs in the first two sessions due to unforeseen time constraints. These sessions, however, did not differ from the following eight in any other regard. In particular, participants in every session signed an identical consent form, which allowed for the possibility of being photographed.

<sup>47</sup>5 - strikingly handsome or beautiful, 2 - above average attractiveness, 3 - average, 2 - plain, below average in attractiveness, 1 - homely, far below average in attractiveness. As in [Biddle and Hamermesh \(1998\)](#), subjects were told to "imagine how the person would look under ordinary circumstances" if they saw an "unflattering facial expression."

Additional controls	Personality of worker	Personality of manager	Sessions	<i>N</i>
None	$p = 0.1543$	$p = 0.2137$	10	156
Output	$p = 0.0630$	$p = 0.6499$	10	156
Physical traits	$p = 0.0054$	$p = 0.6014$	8	123
Output and physical traits	$p = 0.0031$	$p = 0.6393$	8	123

Additional controls	Personality of worker	Personality of manager	Sessions	<i>N</i>
None	$p = 0.0217$	$p = 0.1814$	10	156
Physical traits	$p = 0.0040$	$p = 0.2792$	8	123

Table 19: Correlations between personality and the worker’s effort. Session and half fixed effects are included. The personality traits of the worker remain jointly significant in all but one specification.

	Income	Inequality	Income	Inequality
Agreeableness	-1.304 (1.098)	-2.260*** (0.858)	-1.186** (0.566)	-2.387** (1.197)
Output			1.322*** (0.461)	-1.430 (0.940)
F-statistic (first stage, agreeableness)	24.04	24.04	23.52	23.52
F-statistic (first stage, output)			14.06	14.06
F-statistic (second stage)	1.177	5.791	4.901	2.869
Underidentification test	0.0217	0.0217	0.0240	0.0240
Observations	154	154	154	154

Session-clustered standard errors in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ , \*\*\*\*  $p < 0.001$

Table 20: The effect of the worker’s agreeableness on her bargaining power, with the worker’s rating of the (other) manager used as an instrument. Session half fixed effects are included.

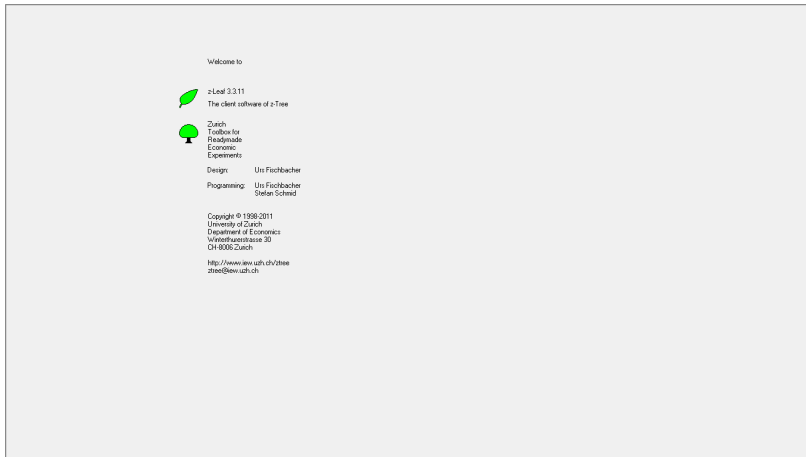
## 4.6 Instructions

In this game, there is a team with Person B and Person A. The role of Person B is

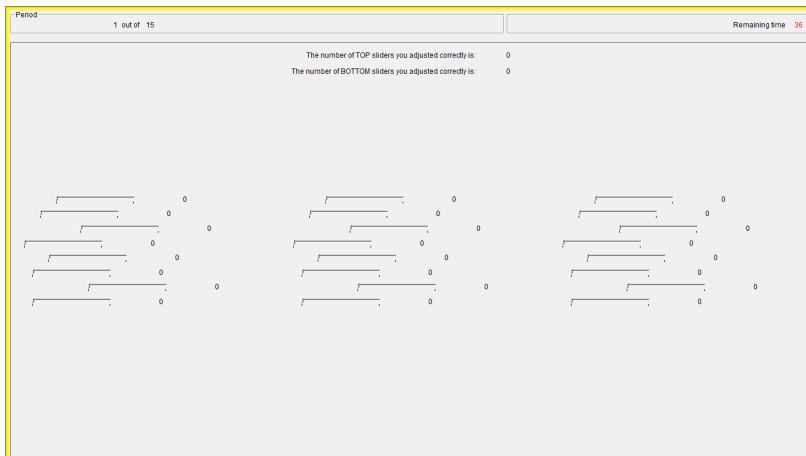
1. to give recommendations to Person A, and
2. to decide how prize money (described below) should be distributed between the team members.

The game will last 15 periods OR until Person B runs out of money.

**In the first period of the game, there is no recommendation.** When the game starts (Period 1 out of 15), Person B will see a screen like this



and Person A will see a screen like this



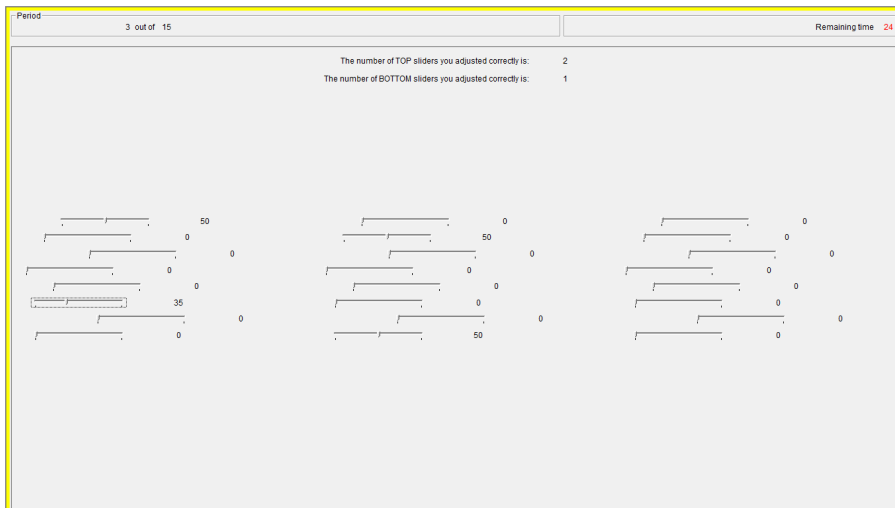


Person A will see 24 sliders on their screen. He/she has **40 seconds** to adjust the sliders.

TOP sliders are sliders in the first four rows. BOTTOM sliders are sliders in the last four rows.

Therefore, Person A will have 12 TOP sliders, and 12 BOTTOM sliders.

Adjusting a slider correctly means adjusting it to position 50. For example, in the screen grab below, one slider in the **first** row has been adjusted to position 50 (correct, TOP), one slider in the **second** row has been adjusted to position 50 (correct, TOP), one slider in the **sixth** row has been adjusted to position 35 (incorrect, BOTTOM), and one slider in the **eighth** row has been adjusted to position 50 (correct, BOTTOM).



Notice that a message in the top part of the screen is informing Person A that two TOP sliders have been adjusted correctly (in the **first** and **second** rows) and one BOTTOM slider has been adjusted correctly (the one in the **eighth** row). The only other slider that has been adjusted – the one in the **sixth** row – has not been adjusted correctly.

# How money is earned

If you are Person A, you earn one penny for each slider NOT at position 50. This money is yours to keep; Person B cannot take it away.

Therefore, in the example above, where three sliders are at 50, Person A will get 21 cents if they keep all sliders at their current positions.

Why would Person A want to adjust sliders at all? One and only one of the 24 sliders contains a prize.

**BUT YOU DON'T GET THE PRIZE MONEY AUTOMATICALLY.** Aside from the pennies Person A gets for unadjusted sliders, Person B is completely in control of the payments received by Person A.

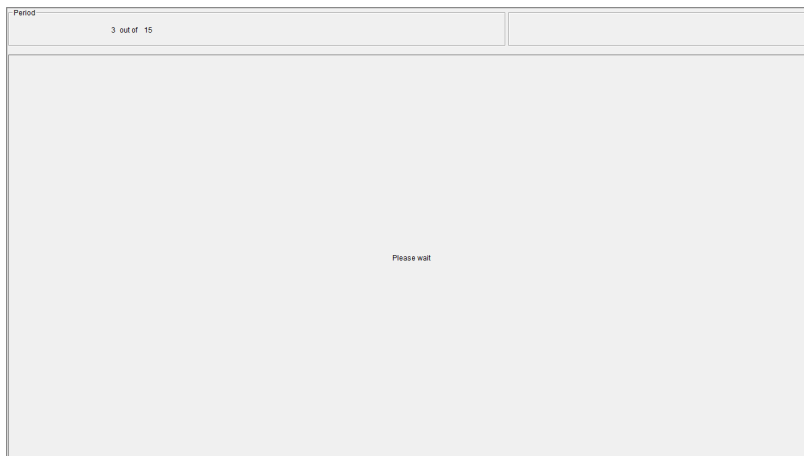
# Where is the prize?

The prize can either be behind a TOP or a BOTTOM slider.

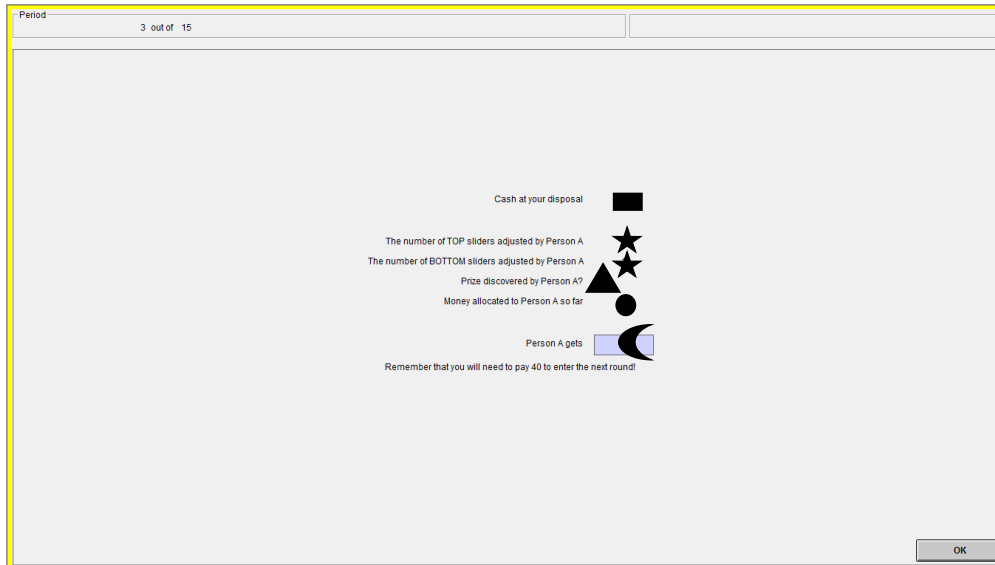
Whether the prize is behind a TOP/BOTTOM slider in the next round only depends on where the prize was in this round.

Person A will never know where the prize is. At the end of every round, Person B will see whether or not the prize was discovered. He/she will use this information to make recommendations to Person A.

After the 40 seconds given to Person A to adjust their sliders run out, Person A will see a screen like this



and Person B will see a screen like this



Now, Person B has to decide how much he/she wants to pay Person A. Person B has unlimited time to make this decision.

In the beginning of the experiment, Person B starts out with 500 cents.

**40 cents are subtracted from Person B's earnings at the beginning of every period.**

Therefore, as soon as Period 1 starts, 40 cents are subtracted from 500, leaving person B with 460.

**If a prize is discovered, Person B gets 400 cents added to their total funds.**

Therefore, you are Person B, and Person A discovered the prize, your available funds (or "Cash at your disposal") at the end of Period 1 will be 860. You will see this number in of the rectangle (■) in the screen grab above.

In place of the triangle (▲), you will see the word "YES" or "NO." YES means that Person A discovered the prize. NO means that Person A did not discover the prize.

Behind the stars (★) you will find information about how many TOP and BOTTOM sliders Person A adjusted.

Person B has to decide how much to pay or fine Person A. This number is entered behind the moon symbol (☾).

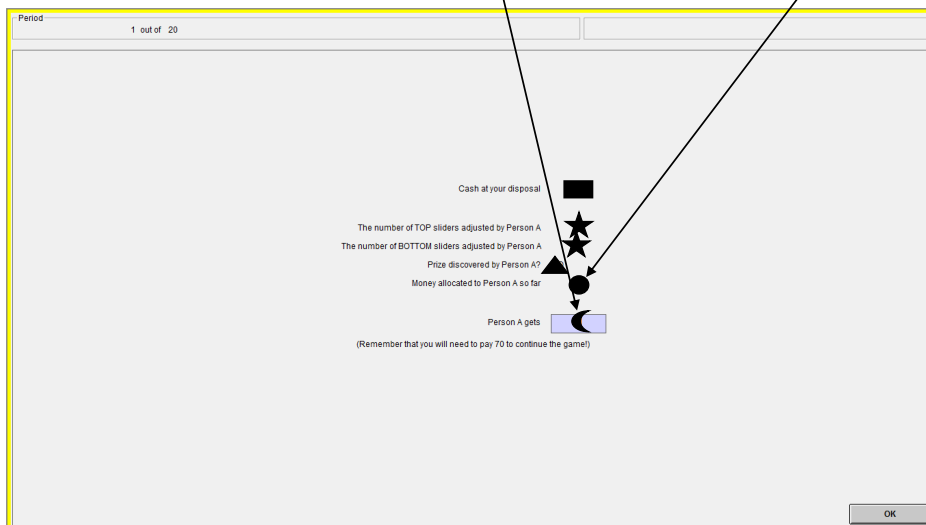
Behind the circle (●) is information about how much Person B paid Person A so far. Whatever is entered behind the moon gets added to the number behind the circle.

# Paying (or fining) Person A

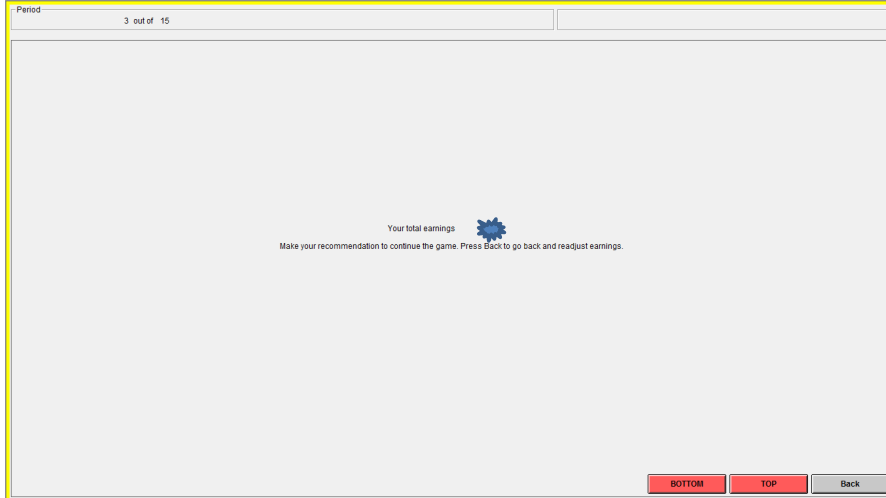
There are a couple of restrictions on how Person B can pay (fine) Person A.

1. **Person B has to make sure that they don't go over the cash at their disposal.** Therefore, if Person B has 1700 cents available, they have to pay Person A no more than 1700.
2. **Person B cannot take more money from Person A than what Person A has been paid so far.**

In other words, whatever is entered [here with a minus sign](#) cannot exceed the number [here](#)

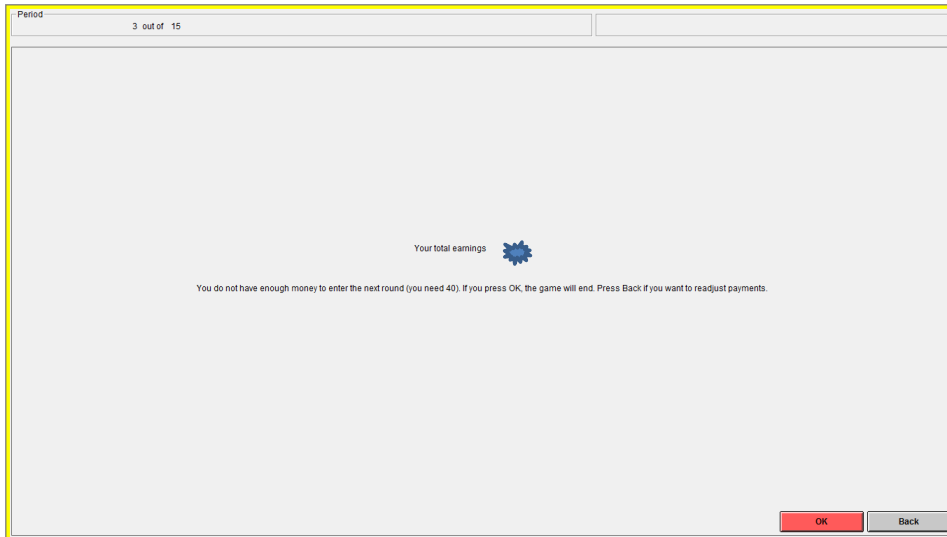


After Person B decides how to pay Person A, **if and only if he/she is left with more than 40**, Person B will see a screen like this



At this point, Person B has to make a recommendation to Person A about which sliders to move.

If Person B has been left with less than 40 after paying Person A, his/her screen will look like this:



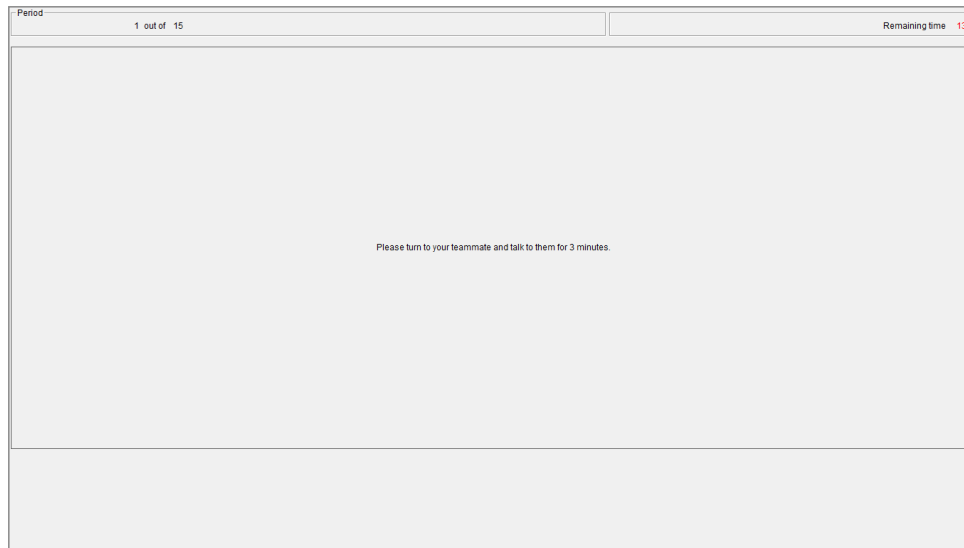
Pressing OK at this point finishes the experiment.

Assuming Person B has been left with more than 40 and the experiment continues, both team members will see their earnings displayed on the screen.

Next, Person A will again see their sliders, and Person B will have to wait 40 seconds before paying Person A and making the next recommendation.

# Communication

At the end of Periods 5 and 10, instead of continuing to the next round, all team members will see the following on their screen.



At this point, Person A and Person B will have three minutes to discuss the game.

**The purpose of these discussions is for Person B to understand how he/she wants to adjust the way he/she has been paying each Person A.**

**You can discuss anything related to the game at this time. As Person B, you share your thoughts and concerns about what Person A is doing. As Person A, you can share your thoughts and concerns about what Person B is doing.**

**You can discuss anything related to the game.**

After the three minutes expire, the experiment will resume.

# Payment

You will be paid **privately**.

Person B will not see how much Person A made. (Although, if he/she keeps count of sliders moved and payments received in the course of the experiment, he/she could calculate this information.)

Person A will not know how much money Person B earned.



## 4.7 Personality Questionnaire

Here are a number of characteristics that may or may not describe you. For example, do you agree that you seldom feel blue? Please fill in the number that best indicates the extent to which you agree or disagree with each statement listed below. Be as honest as possible, but rely on your initial feeling and do not think too much about each item.

Use the following scale:

1 ----- 2 ----- 3 ----- 4 ----- 5  
**Strongly**                      **Neither Agree**                      **Strongly**  
**Disagree**                      **Nor Disagree**                      **Agree**

- |                                                       |                                                   |
|-------------------------------------------------------|---------------------------------------------------|
| 1. ___ Seldom feel blue.                              | 27. ___ Hate to seem pushy.                       |
| 2. ___ Am not interested in other people's problems.  | 28. ___ Keep things tidy.                         |
| 3. ___ Carry out my plans.                            | 29. ___ Lack the talent for influencing people.   |
| 4. ___ Make friends easily.                           | 30. ___ Love to reflect on things.                |
| 5. ___ Am quick to understand things.                 | 31. ___ Feel threatened easily.                   |
| 6. ___ Get angry easily.                              | 32. ___ Can't be bothered with other's needs.     |
| 7. ___ Respect authority.                             | 33. ___ Mess things up.                           |
| 8. ___ Leave my belongings around.                    | 34. ___ Reveal little about myself.               |
| 9. ___ Take charge.                                   | 35. ___ Like to solve complex problems.           |
| 10. ___ Enjoy the beauty of nature.                   | 36. ___ Keep my emotions under control.           |
| 11. ___ Am filled with doubts about things.           | 37. ___ Take advantage of others.                 |
| 12. ___ Feel others' emotions.                        | 38. ___ Follow a schedule.                        |
| 13. ___ Waste my time.                                | 39. ___ Know how to captivate people.             |
| 14. ___ Am hard to get to know.                       | 40. ___ Get deeply immersed in music.             |
| 15. ___ Have difficulty understanding abstract ideas. | 41. ___ Rarely feel depressed.                    |
| 16. ___ Rarely get irritated.                         | 42. ___ Sympathize with others' feelings.         |
| 17. ___ Believe that I am better than others.         | 43. ___ Finish what I start.                      |
| 18. ___ Like order.                                   | 44. ___ Warm up quickly to others.                |
| 19. ___ Have a strong personality.                    | 45. ___ Avoid philosophical discussions.          |
| 20. ___ Believe in the importance of art.             | 46. ___ Change my mood a lot.                     |
| 21. ___ Feel comfortable with myself.                 | 47. ___ Avoid imposing my will on others.         |
| 22. ___ Inquire about others' well-being.             | 48. ___ Am not bothered by messy people.          |
| 23. ___ Find it difficult to get down to work.        | 49. ___ Wait for others to lead the way.          |
| 24. ___ Keep others at a distance.                    | 50. ___ Do not like poetry.                       |
| 25. ___ Can handle a lot of information.              | 51. ___ Worry about things.                       |
| 26. ___ Get upset easily.                             | 52. ___ Am indifferent to the feelings of others. |

53. \_\_\_ Don't put my mind on the task at hand.
54. \_\_\_ Rarely get caught up in the excitement.
55. \_\_\_ Avoid difficult reading material.
56. \_\_\_ Rarely lose my composure.
57. \_\_\_ Rarely put people under pressure.
58. \_\_\_ Want everything to be "just right."
59. \_\_\_ See myself as a good leader.
60. \_\_\_ Seldom notice the emotional aspects of paintings and pictures.
61. \_\_\_ Am easily discouraged.
62. \_\_\_ Take no time for others.
63. \_\_\_ Get things done quickly.
64. \_\_\_ Am not a very enthusiastic person.
65. \_\_\_ Have a rich vocabulary.
66. \_\_\_ Am a person whose moods go up and down easily.
67. \_\_\_ Insult people.
68. \_\_\_ Am not bothered by disorder.
69. \_\_\_ Can talk others into doing things.
70. \_\_\_ Need a creative outlet.
71. \_\_\_ Am not embarrassed easily.
72. \_\_\_ Take an interest in other people's lives.
73. \_\_\_ Always know what I am doing.
74. \_\_\_ Show my feelings when I'm happy.
75. \_\_\_ Think quickly.
76. \_\_\_ Am not easily annoyed.
77. \_\_\_ Seek conflict.
78. \_\_\_ Dislike routine.
79. \_\_\_ Hold back my opinions.
80. \_\_\_ Seldom get lost in thought.
81. \_\_\_ Become overwhelmed by events.
82. \_\_\_ Don't have a soft side.
83. \_\_\_ Postpone decisions.
84. \_\_\_ Have a lot of fun.
85. \_\_\_ Learn things slowly.
86. \_\_\_ Get easily agitated.
87. \_\_\_ Love a good fight.
88. \_\_\_ See that rules are observed.
89. \_\_\_ Am the first to act.
90. \_\_\_ Seldom daydream.
91. \_\_\_ Am afraid of many things.
92. \_\_\_ Like to do things for others.
93. \_\_\_ Am easily distracted.
94. \_\_\_ Laugh a lot.
95. \_\_\_ Formulate ideas clearly.
96. \_\_\_ Can be stirred up easily.
97. \_\_\_ Am out for my own personal gain.
98. \_\_\_ Want every detail taken care of.
99. \_\_\_ Do not have an assertive personality.
100. \_\_\_ See beauty in things that others might not notice.

Use the following scale:

1 ----- 2 ----- 3 ----- 4 ----- 5  
**Strongly**                      **Neither Agree**                      **Strongly**  
**Disagree**                      **Nor Disagree**                      **Agree**

BFAS Scoring Key:

*Neuroticism*

Withdrawal: 1R, 11, 21R, 31, 41R, 51, 61, 71R, 81, 91

Volatility: 6, 16R, 26, 36R, 46, 56R, 66, 76R, 86, 96

*Agreeableness*

Compassion: 2R, 12, 22, 32R, 42, 52R, 62R, 72, 82R, 92

Politeness: 7, 17R, 27, 37R, 47, 57, 67R, 77R, 87R, 97R

*Conscientiousness*

Industriousness: 3, 13R, 23R, 33R, 43, 53R, 63, 73, 83R, 93R

Orderliness: 8R, 18, 28, 38, 48R, 58, 68R, 78R, 88, 98

*Extraversion*

Enthusiasm: 4, 14R, 24R, 34R, 44, 54R, 64R, 74, 84, 94

Assertiveness: 9, 19, 29R, 39, 49R, 59, 69, 79R, 89, 99R

*Openness/Intellect*

Intellect: 5, 15R, 25, 35, 45R, 55R, 65, 75, 85R, 95

Openness: 10, 20, 30, 40, 50R, 60R, 70, 80R, 90R, 100

Reverse response scores for items followed by "R" (i.e. 1=5, 2=4, 4=2, 5=1). To compute scale scores, average completed items within each scale. To compute Big Five scores, average scores for the two aspects within each domain.

Reference:

DeYoung, C. G., Quilty, L. C., & Peterson, J. B. (2007). Between facets and domains: 10 Aspects of the Big Five. *Journal of Personality and Social Psychology, 93*, 880-896.

Contact Colin DeYoung (cdeyoung@umn.edu) for additional information.

## 4.8 Proofs

*Proof of Proposition 1.* The proof follows a basic argument by Fudenberg et al. (1994). For a contradiction, assume that  $\gamma > 20$ . Choose  $v \in E(\delta)$  such that  $v_1 + v_2 = \gamma$ . Player  $i$ 's utility is given by

$$v_i = (1 - \delta)u_i + \delta(pw_i^+ + qw_i^-),$$

where  $w_i^+$  and  $w_i^-$  denote the continuation payoffs of player  $i$  after a good and a bad signal, respectively,  $u_i$  is player  $i$ 's expected utility today and  $p$  and  $q$  are, respectively, the probabilities of a good and a bad signal today. Since  $\gamma > 20$ , it must be the case that the probability of both players cooperating is greater than zero after some history. Let  $\mu_j$ , where  $j \neq i$ , denote player  $j$ 's probability of defection. It will be incentive compatible for player  $i$  to cooperate if

$$\begin{aligned} (1 - \delta)(1 - \mu_j)15 + \delta[((1 - \mu_j)p_2 + \mu_j p_1)w_i^+ + ((1 - \mu_j)q_2 + \mu_j q_1)w_i^-] \\ \geq \\ (1 - \delta)((1 - \mu_j)20 + 2\mu_j) + \delta[((1 - \mu_j)p_1 + \mu_j p_0)w_i^+ + ((1 - \mu_j)q_1 + \mu_j q_0)w_i^-]. \end{aligned}$$

Write  $\Delta p = (1 - \mu_j)(p_2 - p_1) + \mu_j(p_1 - p_0)$ . Rearranging yields:

$$(1 - \delta)[5(1 - \mu_j) + 2\mu_j] \leq \delta \Delta p (w_i^+ - w_i^-),$$

that is, current utility gains from deviating are outweighed future losses in continuation payoffs. Since  $p_2 - p_1 = 1/4$  and  $p_1 = p_0$ , this inequality yields the following upper bound on  $w_i^-$ :

$$w_i^- \leq w_i^+ - \frac{1 - \delta}{\delta} \frac{5(1 - \mu_j) + 2\mu_j}{\Delta p} = w_i^+ - \frac{1 - \delta}{\delta} \frac{5(1 - \mu_j) + 2\mu_j}{(1 - \mu_j)/4} \leq w_i^+ - 20 \frac{1 - \delta}{\delta}$$

Substituting into the previous expression for  $v_i$ ,

$$v_i \leq (1 - \delta)u_i + \delta \left[ w_i^+ - 20q \frac{1 - \delta}{\delta} \right].$$

Therefore,

$$v_1 + v_2 \leq (1 - \delta) [30 - 40q] + \delta \gamma.$$

Since  $q \geq 0.25$  and  $v_1 + v_2 = \gamma$  by hypothesis, it follows that  $\gamma \leq 20$ , as claimed.

Finally, a proof that the bound remains in public communication equilibrium can be found in Rahman (2013a, Lemma 1).  $\square$

*Proof of Lemma 1.* Assume that (3.2) holds. If a player chooses to deviate for  $\tau$  periods, the utility gained from such a deviation is clearly bounded above by  $(1 - \delta)5\tau$ ,

since this bound ignores discounting of future deviation gains. In other words, deviation gains are linear in the number of deviations. On the other hand, punishment costs grow exponentially in the number of deviations. Indeed, the opportunity cost of punishment remains  $\delta^T \alpha(v - 2)$ , but the change in punishment probability from  $\tau$  deviations becomes

$$q_2^{T-\tau}(q_1^\tau - q_2^\tau) = q_2^T \left[ \left( \frac{q_1}{q_2} \right)^\tau - 1 \right],$$

which, since  $q_1 > q_2$ , clearly grows exponentially with  $\tau$ . Now, by the Binomial Theorem,  $(q_1/q_2)^\tau - 1 \geq \tau[(q_1/q_2) - 1] = \tau$ , so the change in punishment probability is bounded below by  $q_2^T \tau$ . Therefore, the following inequality implies that  $\tau$  deviations are discouraged:

$$(1 - \delta)5\tau \leq \delta^T q_2^T \tau \alpha(v - 2).$$

But this is just (3.2). The claim now follows because  $\tau$  was arbitrary.  $\square$

*Proof of Proposition 2.* Fix  $T \in \mathbb{N}$ . As  $\delta \rightarrow 1$ , the left-hand side of (3.3) tends to  $1/T$  by l'Hopital's rule, which is less than or equal to 1. Hence, there exists  $\delta < 1$  sufficiently large that (3.3) holds, so the candidate equilibrium strategies above are indeed an equilibrium. Finally, by l'Hopital's rule, it follows that

$$v \rightarrow 15 - 5\frac{1}{T} \quad \text{as} \quad \delta \rightarrow 1.$$

Finally, it is now clear that  $v \rightarrow 15$  as  $T \rightarrow \infty$ , as claimed.  $\square$

## 4.9 Instructions to treatment NC

### Instructions

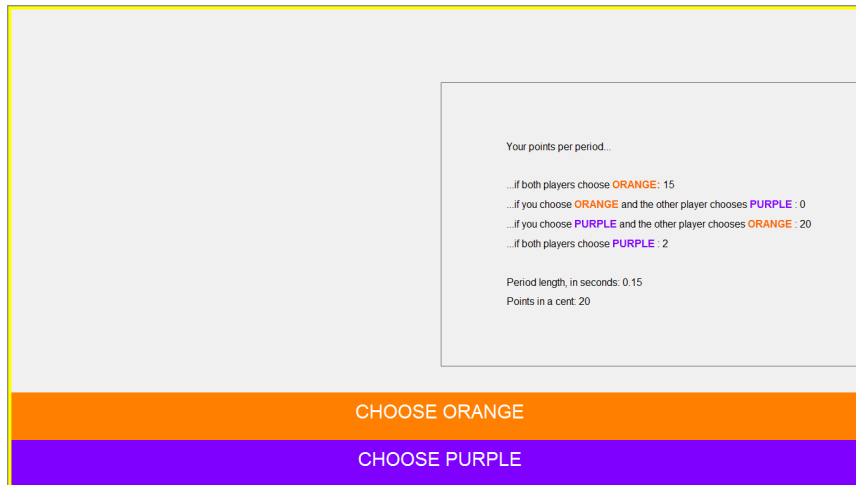
---

Welcome and thank you for participating in the following experiment on strategic decision-making. This document explains what the experiment will entail. First, notice that **your show-up fee that is paid separately and is not affected in any way by the outcome of the experiment.**

#### Timing and Payoffs

You will be randomly and anonymously matched to another subject several times. Each time you are matched to a subject will be called a “match.” During each match, you will interact with the other subject to whom you are matched through a computer program as described below. You will have the opportunity to earn money depending on the decisions made by you and this other subject.

In the beginning of a match, you will see on the bottom of your computer screen an orange button and a purple button. At any time, you will have the choice of selecting either color by clicking on the corresponding button with the computer’s cursor, using your mouse. The image below shows what your computer screen will look like before you make your first choice.



*Initial screen prior to making first choice*

Once you and the other subject have selected an initial color, **periods will start elapsing with a duration of 0.15 seconds per period.** You can change your selection at any time and as often as you want by selecting the corresponding button; so can the other subject to whom you are matched. The computer program will register your changes every period. Thus, in 15 seconds the program will register one hundred choices. **Unless and until you change your choice, your assumed choice for a given period will be the last choice you made.** For example, if you select orange by clicking on the orange rectangle on the screen with your mouse and then change to purple 5 seconds later by clicking on the purple rectangle with your mouse then the computer will register that you chose orange during the time between clicking orange and purple.

Your monetary payoff at the end of the experiment depends both on your color choices and those of the other subjects to whom you were matched. Every period, if you chose orange and the other subject chose orange, too, then you will each earn 15 points. If you chose orange and the other subject chose purple then you will earn zero points and the other subject will earn 20 points. If you chose purple and the other subject chose orange then you will earn 20 points and the other subject will earn zero points. Finally, if both you and the other subject chose purple then you will each earn 2 points. Your final payoff is the accumulation of all your points in all your matches. **Points will be exchanged for money at the rate of forty (40) points per cent, or 1000 points per 25 cents.** The table below summarizes this information.

		Other's choice	
		Orange	Purple
Your choice	Orange	15 points for you 15 points for other	0 points for you 20 points for other
	Purple	20 points for you 0 points for other	2 points for you 2 points for other

*Average points per period depending on each subject's choices*

To illustrate, consider the following example. If you and the other subject to whom you are matched both chose orange for 100 periods, then you would earn  $15 \times 100 = 1,500$  points, translating into  $1,500 \times 1/40 = 37.5$  cents. If you chose purple and the other subject chose orange in every period, you would earn a total of  $20 \times 100 = 2,000$  points, which would translate into  $2,000 \times 1/40 = 50$  cents. If you both chose purple, then you would earn  $2 \times 100 = 200$  points, translating into  $200 \times 1/40 = 5$  cents.

The number of periods in a match is selected as follows. Every period, a random process determines whether the match continues on to the next period. The continuation probability is held constant, so that the average duration of a match is 700 periods. Because termination is random, some matches will last longer than 700 periods and others will last less than that. As soon as a match ends, every subject will be randomly and anonymously re-matched with another subject. You will be re-matched several times. Your final payoff will consist of the accumulation of your payoffs across all matches.

## Information

Throughout a match, you will observe neither your payoff, nor the other subject's payoff, nor the other subject's choices. Similarly, the other subject will observe neither his or her payoff, nor your payoff, nor your choices. You and the other subject will observe the outcome of a random **signal process**, graphically depicted on the left-hand side of your computer screen. The graph of the process will depend on your color selection, the other subject's selection, and an element of randomness, as follows.

Every period, the value of the signal process will either increase or decrease by one unit. If you and the other subject both chose orange, the value of the process will increase with 75% probability and decrease with 25% probability. Otherwise, if one or both of you chose purple then the process will increase and decrease with 50% probability.

The process will be displayed in real time, in blocks of 100 periods. On the top-right region of the screen there will be displayed the fraction of periods during which you chose orange in the current block as well as the **position** of the process, defined as the number of time it actually increased minus the number of times it actually decreased in the current block. If you and the other subject always chose orange, then, at the end of a block, the process will reach a position of around 50 on average. If one or both subjects always chose purple then, at the end of a block, the process will reach a position of around 0 on average. However, this score fluctuates randomly, and can in principle end up far away from these values. At the end of each block, a red "continue" button will appear at the bottom of the screen. You may press the button when you are ready to move on to the next block. There will be two practice blocks at the start to gain familiarity with the process.

To illustrate, see the figures below with possible paths of your earnings over time when you and the other subject make different color choices. Figure 1 below depicts possible paths of the signal process during two consecutive blocks given that both you and the other subject chose orange. The process starts at zero at the beginning of every block. The horizontal graph lines count 20 units of the process increasing or decreasing and the ticker line in the middle denotes the starting point of the signal. In this instance, the signal position exhibited a net rise of 50 units in the first block, followed by a rise of 46.



Figure 1: Possible path of the signal if both of you chose orange throughout

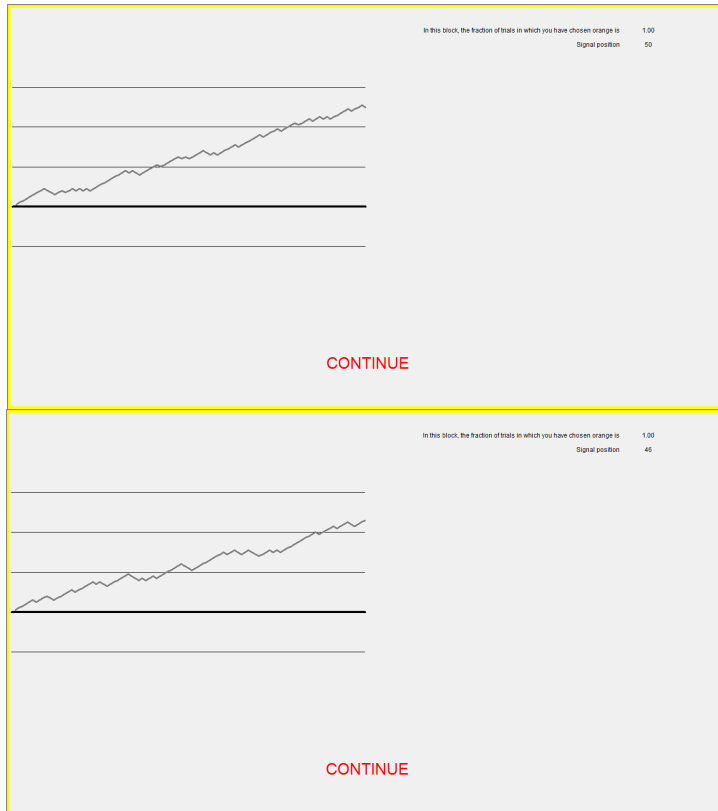


Figure 2: Possible path of the signal if one or both of you chose purple throughout

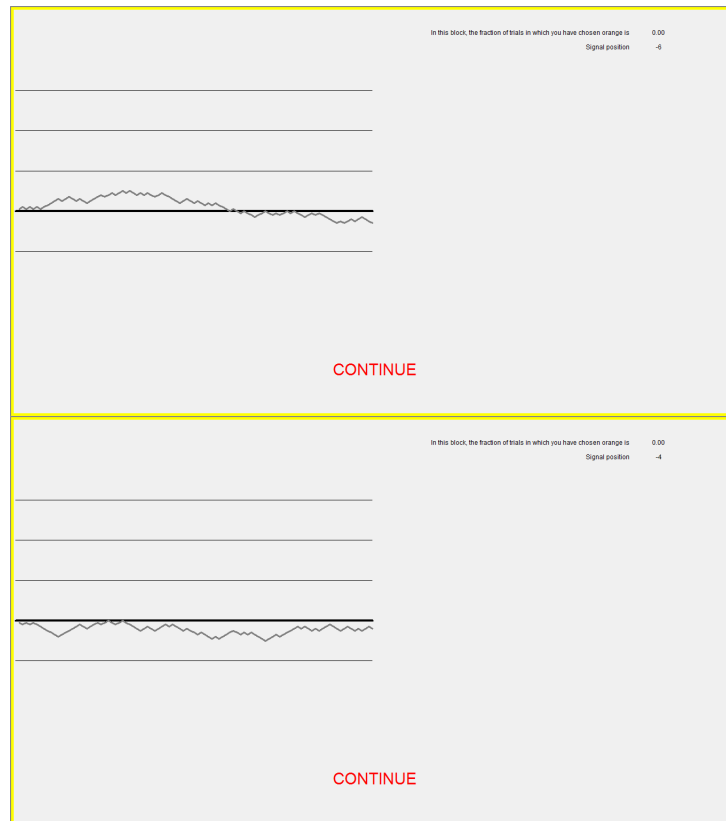
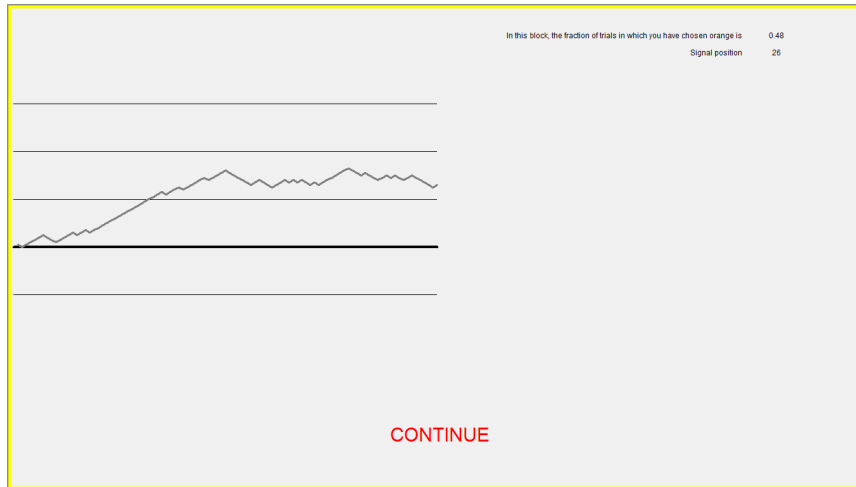


Figure 2 above depicts possible paths of the signal process for two consecutive blocks given that you chose purple and the other subject orange. In this instance, the signal position dropped 6 units in the first block, followed by a drop of 4 units in the next block.

Figure 3 below depicts the possible path of the signal process during a single block given that the other subject chose orange throughout the block and you switched from orange to purple halfway through the block.

*Figure 3: Possible path of the signal if you switch from orange to purple after 48 periods*



## Communication

In addition to observing the signal process, you will be able to send messages to and receive messages from the other subject. At the beginning of every block, you will be able to tell the other subject the percentage of time you plan to choose orange in the next block. You will be able to tell the other subject a plan for choosing orange some percentage of time if the signal position is above or below some number of your choosing. Once you have entered and submitted your answers, your message will be sent to the other subject and you will receive the other subjects' message. You will see on the right-hand side of your screen both your most recent message and the other subject's most recent message throughout the next block. At the end of every block you will observe your most recent message as well as the other subject's most recent message while completing your next message for the subsequent block. You will have the option of submitting the same message as before or submitting a different message. Below are some screenshots to illustrate.

Figure 4: Screenshot of initial message screen at the beginning of the first block

Your points per period...

- ...if both players choose **ORANGE**: 15
- ...if you choose **ORANGE** and the other player chooses **PURPLE**: 0
- ...if you choose **PURPLE** and the other player chooses **ORANGE**: 20
- ...if both players choose **PURPLE**: 2

Period length, in seconds: 0.15  
Points in a cent: 20

This block, I will choose **ORANGE** this percentage of the time

If this block's signal position is  ABOVE  
 BELOW

the number

I will respond by choosing **ORANGE** this percentage of the time in the following block

Otherwise, I will choose orange this percentage of the time

Figure 5: Screenshot of messages sent and received at the beginning of the first block

Remaining time: 27

The other player said that he/she will choose **ORANGE** this percentage of the time in this block: 2

The other player also reported that if the signal position is **BELOW**  
the number: 2  
he/she will respond by choosing **ORANGE** this percentage of the time in the next block: 2  
and that otherwise, he/she will choose **ORANGE** this percentage of the time: 2

You said that you will choose **ORANGE** this percentage of the time in this block: 1

You also reported that if the signal position is **ABOVE**  
the number: 1  
you will respond by choosing **ORANGE** this percentage of the time in the next block: 1  
and that otherwise, you will choose **ORANGE** this percentage of the time: 1

Figure 6: Screenshot of a subject at the end of a block



Figure 6: Screenshot of a subject at the end of a block if changing message



## Ground Rules

Please wear the headphones provided throughout the experiment, except when instructed to do so by the experimenters. We also ask that you disconnect your cellphones throughout the duration of the experiment.