

Promises and Punishment^{*}

Martin Dufwenberg[†] Flora Li[‡] Alec Smith[§]

March 25, 2025

Abstract

We study the effect of communication on beliefs, behavior, and efficiency in the context of hold-up problems with a punishment option. We apply a novel behavioral motivation, frustration-dependent anger, that links unmet payoff expectations with the willingness to forgo material payoffs to punish others, and we conjecture that communication works through this mechanism to raise expectations about the likelihood of belief-dependent costly punishment and to increase trust, cooperation, and efficiency. In an experiment we allow communication in the form of a single pre-play message. We measure beliefs and our design permits the observation of promises and deception. The results are consistent with the theory that costly punishment results from belief-dependent anger and frustration. Promises drive the effect of communication on beliefs and broken promises lead to higher rates of costly punishment.

Keywords: Communication, Hold-up, Frustration and anger, Promises, Punishment, Psychological game theory

^{*}We are particularly grateful to Pierpaolo Battigalli, Gary Charness, Uri Gneezy, and many seminar and conference participants for their helpful comments. The Center for Peace Study and Violence Prevention at Virginia Tech provided financial support.

[†]University of Arizona; University of Gothenburg; CESifo. Email: martind@eller.arizona.edu.

[‡]Corresponding author. Guangdong Institute of Intelligence Science and Technology. Email: lizhuncheng@gdiist.cn.

[§]Department of Economics, Virginia Tech. Email: alecsmith@vt.edu.

1 Introduction

Communication can foster trust and cooperation. A recent literature explores why people keep their promises, focusing on the motivation of the promisor.¹ We explore a new and complementary explanation, whereby it is the promisee that is affected. If a promise is broken this induces dashed hopes and frustration, which triggers anger and aggression.² If anticipated, this creates an incentive for promisors not to renege.

We develop this idea for environments which augment a simple trust game with a punishment stage. The resulting structures may be viewed as particular forms of hold-up problems, where relationship-specific investments and incomplete contracts expose one party to opportunistic renegotiation, potentially resulting in underinvestment.³ We explore such settings both theoretically and experimentally:

First, we apply the model of frustration and anger from Battigalli, Dufwenberg, and Smith (2015; 2019) (BDS) to a three-stage hold-up problem, allowing us to examine the impact of communication in general and promises in particular. The basic ideas are: 1) decision-makers experience anger when they are frustrated; 2) frustration results when material payoffs are less than expected; and 3) anger leads to aggression and the urge to retaliate; 4) promises may enhance the just-mentioned effects, by shaping expectations such that promise-keeping is expected; 5) if these effects are anticipated, trust and cooperation ensues. The approach requires a formulation of utility where a player's preferences depend both on material payoffs and on beliefs about his own and others' behavior.⁴ Messages become relevant to the extent that they influence expectations about payoffs, thus linking communication, beliefs, and the willingness to forgo material payoffs to punish others.

Second, we design an experiment to test the predictions of the theory.⁵ We allow pre-play communication as a treatment in order to study whether promises are sent and whether their

¹Charness and Dufwenberg (2006) develop a theoretical argument based on “guilt aversion” and Vanberg (2008) similarly explores “a preference for keeping one’s word.” Some of the large subsequent literature is surveyed by Cartwright (2019) and Di Bartolomeo et al. (2023).

²Psychologists associate frustration with aggression; see e.g. Dollard et al. (1939); Berkowitz (1989).

³See Williamson (1971); Klein et al. (1978); Grout (1984); Grossman and Hart (1986); Tirole (1986); North and Weingast (1989); and Hart and Moore (1990); compare *e.g.* Ellingsen and Johannesson (2004a), Ellingsen and Johannesson (2004b), Che and Sákovics (2008), and Dufwenberg et al. (2013) who explain how the setup we consider involves the sub-class of hold-up problems with a punishment option.

⁴The approach involves belief-dependent utilities and draws on the framework of psychological game theory (Geanakoplos et al., 1989; Battigalli and Dufwenberg, 2009, 2022).

⁵BDS (2019, pp. 17, 29, 31) discuss previous attempts by economists to address frustration and anger, either theoretically or experimentally. Two of the experimental studies - Persson (2018) and Aina et al. (2020) - relate directly to BDS, although unlike us these authors do not explore issues of communication.

effect on beliefs and behavior is as we predicted. A key contribution of our paper is that, in addition to recording messages and behavior, we elicit the beliefs of both players before messages are sent and after they are received. We measure beliefs about co-player choices, and also, in a novel contribution, about players’ own behavior at subsequent stages of the game. These measures allow us to carefully examine the relationship between communication, beliefs, and behavior. In particular, we measure how promises change beliefs and how expectations about behavior influence the decision to engage in costly punishment.

Ellingsen and Johannesson (2004a,b), who also study communication and hold-up in an experiment, are important precursors to our study. However, since they did not conduct their exercise with BDS’ theory in mind, they did not measure the beliefs which are central to our tests.⁶ Less closely related are several experimental studies of hold-up games that do not focus on the impact of communication. See Yang (2021) for a recent review.

Section 2 presents theory. We describe the games we study, apply BDS’ model of belief-dependent anger, and discuss the extension needed to incorporate the ideas we have regarding the effect of promises on trust, credibility, and costly punishment. Section 3 presents details of the experimental design and implementation, and states hypotheses to be tested. Section 4 reports summary statistics, main results regards hypotheses, and additional observations. In Section 5 we discuss alternative motivations, and Section 6 concludes.

2 Theory

2.1 A hold-up game with costly punishment

We study a class of 2-player, 3-stage games, as shown in Figure 1, where the numbers and variables at end nodes represent monetary payoffs. The game may be interpreted as a mini-trust game with an added (subsequent) punishment option, an ultimatum mini-game with an added (preceding) entry decision, or as a hold-up game where sellers can destroy the proceeds of a relationship-specific investment.⁷ In the first stage, Player 1 can choose *In*

⁶Ellingsen and Johannesson suggest that their data is consistent with Fehr and Schmidt’s (1999) model of inequality aversion combined with a preference for consistency, and that communication serves to change beliefs about co-player types. This interpretation is quite different from the theory that we test. Later on, we address how models of inequity aversion relate to our data.

⁷In general, hold-up may occur in environments with or without the opportunity for punishment or “vengeance” (Dufwenberg et al., 2013). In order to study of the effect of broken promises we focus on a hold-up environment that allows for costly punishment after opportunistic behavior.

to make an investment of her entire endowment of \$5, or *Out* to not invest and walk away with her initial endowment. If Player 1 invests, the endowments of both players double, and Player 2 can then propose how to divide the proceeds. To make the problem simple, Player 2 can propose one of two possible splits. One option is to choose *Share*, which is monetarily favorable (or at least as good as the other option) for Player 1. The other is to choose *Take*, which is (potentially) monetarily favorable for Player 2. If Player 2 *Takes*, Player 1 can then *Reject*, in which case both players receive 0, or *Accept* to settle with a less favorable offer in the third stage. The parameters a and b reflect the payoffs to Player 1 after, respectively, $(In, Share)$ and $(In, Accept), Take)$. We impose the following parameter restrictions: $20 > a \geq 5 \geq b > 0$, and $a \neq b$. When players care only for monetary payoffs and $b < 5$, the unique subgame perfect equilibrium (SPE) is $((Out, Accept); Take)$, which is inefficient; when $b = 5$ and players care only for monetary payoffs, there are two SPEs: $((Out, Accept); Take)$ and $((In, Accept); Take)$.

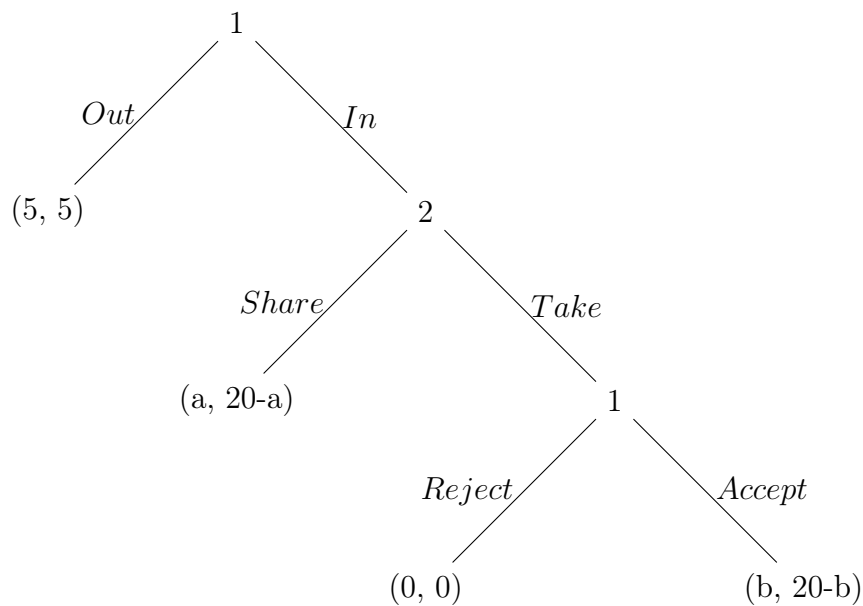


Figure 1. A hold-up game with punishment.

2.2 Frustration and anger

We apply BDS’s frustration and anger model.⁸ In this model, anger is motivated by frustration, and the tendency to hurt others is proportional to frustration, following the frustration-aggression hypothesis from psychology (Dollard et al., 1939; Berkowitz, 1989). In general, one feels frustrated if one’s initial expectation is not met. Frustration is modeled as the gap (if positive) between one’s initial expected payoff and the current best possible outcome. At any history h , Player i ’s frustration is

$$F_i(h; \alpha_i) = \max \left\{ \bar{\pi}_i(h_0) - \max_{a_i \in A_i(h)} \mathbb{E}[\pi_i|h; \alpha_i], 0 \right\}, \quad (1)$$

where $\bar{\pi}_i(h_0) = \mathbb{E}[\pi_i|h_0; \alpha_i]$ denotes Player i ’s expected payoff at the initial history h_0 given his first-order belief α_i about Player j ’s behavior, $a_i \in A_i(h)$ denotes Player i ’s action choice at the history h , so $\max_{a_i \in A_i(h)} \mathbb{E}[\pi_i|h; \alpha_i]$ gives the maximum possible expected payoff available to Player i at the history h .

Player i ’s utility from action a_i at history h is

$$u_i(h, a_i; \alpha_i) = \mathbb{E}[\pi_i|(h, a_i); \alpha_i] - \theta_i F_i(h; \alpha_i) \mathbb{E}[\pi_j|(h, a_i); \alpha_i], \quad (2)$$

where $\theta_i \geq 0$ is Player i ’s sensitivity to anger. A frustrated individual is motivated to hurt the other player, if the cost is low enough. Frustration increases the negative weight placed on the Player j ’s material payoff, and motivates aggression.

In the game forms defined in Figure 1, Player 1 is the party who might get frustrated, so we apply Equations (1) and (2) with $i = 1, j = 2$. Let the probability that Player 1 assigns to choosing *Out* be $p_1 = \alpha_1(Out|h^0) \in [0, 1]$. Let $q_1 \in [0, 1]$ denote the probability that Player 1 assigns to Player 2 choosing *Share* if stage 2 is realized, i.e. $q_1 = \alpha_1(Share|In)$ and let $r_1 = \alpha_1(Reject|In, Take) \in [0, 1]$ denote the probability that Player 1 assigns to choosing *Reject* conditional on the 3rd stage being reached. We can also define analogously a similar belief system (p_2, q_2, r_2) for Player 2. We further assume that beliefs are coherent, so the marginals of the higher order beliefs are equal to the lower order beliefs.

⁸BDS model three versions of belief-dependent frustration and anger: 1) Simple anger (SA), 2) Anger from blaming behavior (ABB), and 3) Anger from blaming intentions (ABI). In the hold-up environment studied here, the predictions of all three models coincide (although the math below reflects the SA-formulation). BDS (2019) focus on two-stage “leader-follower” games; an earlier working paper BDS (2015, Section 6) develops the extension to general multi-stage games.

Next, we derive equilibrium predictions, applying the sequential equilibrium (SE) concept from Battigalli et al. (2019).⁹ With utility as defined in Equation 2 there are two pure-strategy SEs of this game: an efficient one, where Player 1 chooses *In*, Player 2 *Shares*, and Player 1 *Accepts*; and an inefficient one, which coincides with the subgame perfect equilibrium for material-payoff maximizing players.

In any SE, beliefs must be correct in the sense that they are consistent with behavior. This means that the first order belief of player i about what player j will do match player j 's behavior strategy. In addition, in equilibrium the belief systems of both players coincide, so for expedience, we drop the subscripts and generically refer to beliefs p, q , and r . We focus here on SE's involving pure strategies.

If θ_1 is small, then the unique SE coincides with the SPE for players who only care for material payoffs: $((Out, Accept); Take)$, with beliefs $p = 1, q = 0, r = 0$ for both players. Player 1's initial expected material payoff is $5p + a(1 - p)q + b(1 - p)(1 - q)(1 - r) = 5$. With these beliefs, frustration equals $5 - b$ after *Take*. In that case Player 1 compares 0 (the payoff from *Punish*) to $b - \theta_1(5 - b)(20 - b)$ (the payoff from *Accept*), and chooses *Accept* if $\theta_1 < \frac{b}{(5-b)(20-b)}$. We refer to this SE as the "inefficient equilibrium."

If Player 1's sensitivity to anger θ_1 is sufficiently large, this (psychological) game has a unique SE involving the strategy profile $((In, Reject); Share)$ where Player 1 chooses *In*, Player 2 chooses *Share*, and if Player 2 instead chooses *Take* then Player 1 chooses *Reject*. For $((In, Reject); Share)$ to be an SE, the correct beliefs system is $p = 0, q = 1, r = 1$ for both players. Player 1's initial expected material payoff is $5p + a(1 - p)q + b(1 - p)(1 - q)(1 - r) = a$, and at the history $(In, Take)$ Player 1's frustration equals $a - b$. If he gets the move after *Take*, Player 1 then compares the payoff of 0 from choosing *Reject* to the payoff $u_1 = b - \theta_1(a - b)(20 - b)$ from *Accept*. Given equilibrium beliefs, Player 1 will *Reject* if $\theta_1 > \frac{b}{(a-b)(20-b)}$, demonstrating the uniqueness of the efficient equilibrium for large θ_1 . We refer to this as the "efficient equilibrium."

For intermediate values of θ_1 , both the inefficient and efficient equilibrium exit. To see why, recall that Player 1's frustration following *Take* is $5 - b$ in the inefficient SE and $a - b$ in the efficient SE. Since $a \geq 5$, then $5 - b \leq a - b$ (with strict inequality if $a > 5$). Hence, the lowest value of θ_1 that makes it a best response for Player 1 to *Reject* in an efficient SE is

⁹The SE concept was extended to psychological games by Battigalli and Dufwenberg (2009). Battigalli et al. (2019) focus on leader-follower games. The game form in the present study is not a leader-follower game, but the definitions in that paper extend naturally. For a full development, see Battigalli et al. (2015, Section 6).

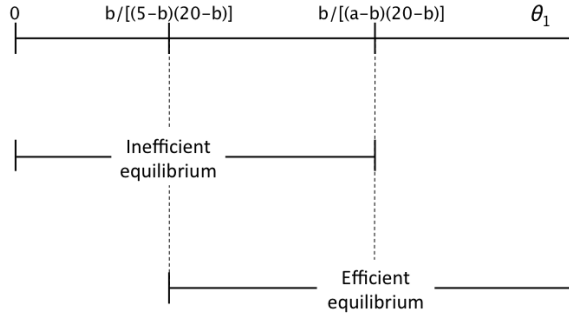


Figure 2. Sequential Equilibria, as a Function of the Anger Sensitivity θ_1 of Player 1.

at least as low as the highest value of θ_1 that make it a best response for Player 1 to *Accept* in an inefficient SE.

The anger sensitivity θ_2 of Player 2 plays no role in the analysis. If Player 2 gets the move, then her maximal payoff is still available and according to the model, she cannot be frustrated ($F_2(\cdot) = 0$) and her behavior is indistinguishable from material payoff maximization. Therefore we focus our analysis on the beliefs and behavior of Player 1.

To summarize: For small values of θ_1 , the unique SE coincides with the inefficient subgame perfect equilibrium for material-payoff maximizers. For large values of θ_1 , the unique SE is efficient: Player 2 *Shares* to avoid being punished, and so Player 1 chooses *In*. For intermediate values of θ_1 there exist two SE in pure strategies, the inefficient one and the efficient one. Figure 2 summarizes how these SE map to the anger sensitivity of Player 1.

2.3 Communication and Promises

If the players are selfish (in particular, if $\theta_1 = 0$, since as noted θ_2 is not relevant) and $b < 5$, the game has a unique backward induction solution: $((Out, Accept); Take)$. The logic of that prediction is not affected by whether or not there is an opportunity for pre-play communication and promises.

By contrast, if $\theta_1 > 0$, promises may plausibly affect behavior and beliefs. We predict that a promise-to-*Share* will (i) increase the likelihood of choices *In*, *Share*, and *Reject*, and (ii) that p_i will decrease while q_i and r_i will increase, for $i = 1, 2$.

A special case is an instance of equilibrium selection, if $\frac{b}{(5-b)(20-b)} < \theta_1 < \frac{b}{(a-b)(20-b)}$ (compare Figure 2). Recall that there are two SE, $((Out, Accept); Take)$ and $((In, Reject); Share)$. Assume that absent communication players play $((Out, Accept); Take)$. Assume that following a promise-to-*Share* players play $((In, Reject); Share)$. In this case a promise-to-*Share* allows the players to coordinate on the Pareto-efficient equilibrium.¹⁰

Predictions (i) and (ii) are not limited to equilibria though. Rather, the idea is that messages feed self-fulfilling chains of beliefs that better choices will be made more frequently. In particular, suppose Player 2 issues a promise-to-*Share*:

- If Player 1 attaches credibility to Player 2’s promise, then q_1 rises.
- With a higher q_1 , Player 1’s frustration following $(In, Take)$ increases (see Equation 1).
- Player 1’s increased frustration makes *Reject* a better choice for Player 1, suggesting an increased frequency of *Reject* as well as higher r_i ’s.
- The increased r_2 makes *Share* a better choice for Player 2, suggesting an increased frequency of *Share* as well as higher q_2 .
- When q_1 and r_1 increase, as long as the increase in q_1 is large enough, this makes *In* a better choice for Player 1, suggesting an increase in p_i .

To sum it up, under the assumption that certain messages influence beliefs, the belief-dependent frustration-anger model implies that, with promises, Player 1 is more likely to trust Player 2 (choose *In*), Player 2 is more likely to keep her promises (to *Share*), and Player 1 is more likely to punish broken promises (by *Rejecting*).

3 Experiment

To study the effect of promises on trust and punishment we implemented a laboratory experiment with the class of games depicted in Figure 1. We employed a within-subject design where subjects played variations of the game over multiple rounds, with fixed roles,

¹⁰See Crawford (2016) for a broad discussion of how there is some empirical support for communication to have such efficiency-enhancing effects, as well as a nuanced critical discussion of the difficulty in establishing clear game-theoretic underpinnings.

paired with anonymous partners with random rematching each round. Each session included a communication and a no-communication block, with the order counterbalanced across sessions.

3.1 Procedures

The experiment was programmed using z-Tree (Fischbacher, 2007) and conducted at the Virginia Tech Economics Laboratory. A sample of the experiment instructions is reproduced in the Appendix. We conducted a total of 11 sessions, with 200 total participants.¹¹ Each session included 14-20 participants with an average of 18.4 per session. Sessions took about 1.75 hours to complete.

At the beginning of each session, participants were randomly assigned to the role of either Player 1 or Player 2, which remained fixed throughout the experiment. Before each round, participants were randomly and anonymously matched with a partner of the opposite role (i.e., we used stranger matching). After the experiment, participants were paid according to the outcome of one randomly selected round. Excluding the show-up fee, participants earned an average of \$12.24.¹²

Each session consisted of 20 rounds separated into two blocks: 10 rounds of communication, and 10 rounds where no communication was allowed. After each round both players were informed of the outcome. We counterbalanced the order of the communication block across sessions, so that in 5 of the 11 sessions the first 10 rounds involved pre-play messages from Player 2 to Player 1, and the no-message block followed; the other 6 sessions experienced the no-message block first. The only restrictions on message sending were that the message had to be less than 140 characters long, and to retain confidentiality, individuals were not allowed to reveal their identity in the message.

In each block, participants played 10 different variations on the game in Figure 1, in random order. The game variations are shown in Table 1, where all the numbers are in US dollars. A change of the parameter b (Payoff from *Accept*) indicates changing the cost of *Reject*, and we vary the cost of *Reject* from 1 to 5. The difference $a - b$ indicates the “*Take* amount”, which takes one of two values in our design: either $a - b = 4$ to indicate a low *Take* amount, or $a - b = 10$ to indicate a high *Take* amount. The payoff splits in Stage 2

¹¹We dropped the data from one additional session that was interrupted by a software malfunction.

¹²After Session 4, we increased the show-up fee from \$5 to \$10 to improve turnout.

and Stage 3 are asymmetric, such that $a \neq 10$, to reduce the saliency of an equal split.

Table 1. Experiment Design – Game Structure.

Game	Value of <i>Share</i> (a)	Cost of <i>Reject</i> (b)	<i>Take</i> Amount ($a - b$)
LT1	5	1	4
LT2	6	2	4
LT3	7	3	4
LT4	8	4	4
LT5	9	5	4
HT1	11	1	10
HT2	12	2	10
HT3	13	3	10
HT4	14	4	10
HT5	15	5	10

It is common in experiment designs to make use of the “strategy method”, where players make conditional choices before the game is played. However, the BDS model implies that players experience zero frustration in this set-up, as frustration can only arise in the course of play. Consistent with our motivation, Aina et al. (2020) demonstrate that frustration and anger are more relevant with the direct response method than with the strategy method, and that costly punishment exhibits greater belief-dependence in sequential decisions. In addition, Brandts and Charness (2011) show that costly punishment is more frequently observed with the direct response method than with the strategy method. Accordingly, we used a direct response design, such that players move sequentially.

3.2 Belief Elicitation

During the experiment we elicited each participant’s probabilistic beliefs about their co-player’s actions, conditional upon future play. We also asked participants to report the probability with which they expect to take an action conditional upon future play in the game. In each round, we measured the first-order beliefs that participants held about their own (in the case of first movers) and their co-players’ behavior. We elicited Player 1’s beliefs regarding the likelihood of choosing *Out* (p_1), Player 1’s conditional first order beliefs of Player 2’s probability of choosing *Share* (q_1), and Player 1’s beliefs regarding the likelihood of choosing *Reject* (r_1) conditional on entering the 3rd stage. We interpret players’ beliefs about their own choices as revealing their *plans*. Our data on plans allow us to gauge players’ would-have-been behavior at nodes that are not actually reached when the game is played out (e.g., if 1 chose *Out*), despite that we did not use the strategy method.

To examine how messages influence beliefs, in the communication treatment we measure Player 1’s beliefs both before and after messages are received. If Player 1 chose *In*, we elicited Player 2’s second order belief about *Share* (q_2) and first order belief about the conditional probability that Player 1 will choose *Reject* (r_2) after Player 2 made a decision on the 2nd stage.

When combining the direct response method with belief elicitation, as we do, there is a potential conflict between incentives for behavior and for reporting truthful beliefs (e.g. Rutström and Wilcox, 2009; Blanco et al., 2010).¹³ In sequential play designs, incentivizing truthful belief reports by *e.g.* a scoring method can create a spillover effect where players have incentives to continue the game in order to receive payment for a reported belief in a future stage, or to choose actions that are consistent with a previously reported belief. The problem is exacerbated when eliciting beliefs about a player’s own future behavior.¹⁴ Trautmann and Kuilen (2015) find that flat fee incentives perform almost as well as more complicated methods for eliciting beliefs such as proper scoring rules. We therefore eschew the use of a scoring rule for payment, instead incentivizing belief reports with a flat fee payment of \$5. In addition, we asked participants to pledge to answer these questions “to the best of my knowledge,” in an attempt to trigger a desire for honest response (see the instructions in the Appendix).

3.3 Hypotheses

Our hypotheses are based upon two main assumptions. First, participants in the experiment are motivated by belief-dependent anger. Second, messages are informative.

With regard to belief-dependent anger, the model implies that unmet expectations regarding material payoffs will increase the disutility from a co-player’s payoff. Hypothesis 1 is motivated by the theoretical assumption that diminished payoff expectations make aggression and costly punishment more attractive (Section 2.2).

Hypothesis 1. *Reject choice frequencies and plans to Reject (r_1) are increasing in Player 1’s belief about the probability of Share (q_1).*

¹³See also Schotter and Trevino (2014) for a review of the methodology of eliciting beliefs.

¹⁴See also the discussion of incentivizing own beliefs in Toussaert (2018), who addresses this issue by eliciting beliefs about a “similar other.” Because we are interested in *own* beliefs as the relevant variable for anger and costly punishment, we also do not employ methods that involve proxies such as similar others or the average belief in the room (as in Charness and Dufwenberg, 2006).

Our model predicts that beliefs and payoffs interact to influence Player 1's behavior (see Section 2.2). The efficient SE $((In, Reject); Share)$ is unique when the anger sensitivity parameter is sufficiently large, such that $\theta_1 > \frac{b}{(a-b)(20-b)}$. This expression shows that when the *Take* amount $(a-b)$ is high, the minimum value of θ_1 that supports the choice of *Reject* in the efficient equilibrium is lower. Similarly, as the cost of *Reject* (b) increases, higher values of θ_1 are necessary to support the choice to *Reject*. These comparative statics motivate the following:

Hypothesis 2. *Player 1's are less likely to Reject (and plan to Reject, r_1) as the cost of Reject (b) increases. Player 1's are more likely to Reject (and plan to Reject) when the Take amount $(a-b)$ is high.*

We next turn to the effect of communication. With reference to the arguments made in Section 2.3, we expect that communication will increase the frequency of cooperative outcomes and lead to greater total material payoffs (which we refer to as efficiency). This is furthermore consistent with a number of studies of communication and cooperation (Charness and Dufwenberg, 2006; Vanberg, 2008; Balliet, 2010), and studies of communication and efficiency (Blume and Ortmann, 2007; Avoyan and Ramos, 2020; Fehr and Sutter, 2019).

Hypothesis 3. *Communication increases the frequency of the cooperative outcome $(In, Share)$ and improves efficiency.*

Motivated by theoretical description in Section 2.3 and the results of Charness and Dufwenberg (2006) and the subsequent literature, we hypothesized that the content of the free-form messages would play an important role in connecting communication with behavior via beliefs. In particular, we predicted that promises would change beliefs and plans in the direction of increased investment, cooperation, and punishment. Hypotheses 4 and 5 connect communication and costly punishment via the effect of communication on beliefs. With regard to message content, we hypothesize:

Hypothesis 4. *Communication influences beliefs via promises, such that promises shift Player 1's reported beliefs and plans in the direction of increased likelihood of In (p_1), $Share$ (q_1), and $Reject$ (r_1); non-promises have no impact on beliefs.*

We predicted that the effect of promises on beliefs would carry through to behavior, through the mechanism of belief-dependent anger as described in Section 2. In particular, an implication of the frustration-anger model is that if promises are believed and then broken, the higher initial expectation of cooperation generates greater frustration and leads to a higher likelihood of rejection in the 3rd stage:

Hypothesis 5. *Broken promises lead to a higher *Reject* rate, and promises lead to a higher cooperation (*Share*) rate relative to non-promises.*

4 Results

We begin our examination of the results by reporting summary statistics on behavior in Section 4.1, and then present our main results following the order of our preconceived hypotheses (Section 4.2). We report additional observations in Section 4.3.

4.1 Data

Our dataset includes choices, elicited beliefs and plans, and messages (when communication was available) from 11 experiment sessions involving 200 participants total who each participant made choices in 20 rounds of game play. Figure 3 shows aggregate results for each stage, pooling participants and games. Participants chose *In* 1,542/2,000 times (77.1%). This is clearly much greater than the upper bound of 20% implied by the SPE for selfish players.¹⁵ Of games that advanced to the 2nd stage, participants chose *Share* 1,068/1,542 times (69.3%). In the third stage, participants selected *Reject* 179/474 times (37%). In both stages 2 and 3, the SPE prediction for selfish players is unique for all 10 games, involving *Take* and *Accept*. Thus, the majority of our observations involve departures from purely selfish behavior.

Behavior. Figure 4 shows the relative frequencies of *Out*, *Share*, and *Reject* choices, arranged by the cost of *Reject*(b) and the *Take* amount.

In Low Take games (LT1-LT5), the relative frequency of *Out* choices declines from 52% (when the cost of *Reject* (b) is 1) to 6.5% (when the cost of *Reject* is 5). In High Take games (HT1-HT5), the relative frequency of *Out* choices is 16.5% when $b = 1$, rising to 27% when $b = 5$).

In all games we observe a decreasing tendency for Player 2's to select *Share* as b increases. Share rates are close to 1 when $b = 1$. When $b = 5$, *Share* is chosen nearly 70% of the time in the Low Take game (LT5) but about 20% of the time in the High Take game (HT5). These

¹⁵This upper bound of 20% is due to Games LT5 and HT5, where Player 1's give up a payoff of 5 to select *Reject*, the same amount that could have been earned by selecting *Out*.

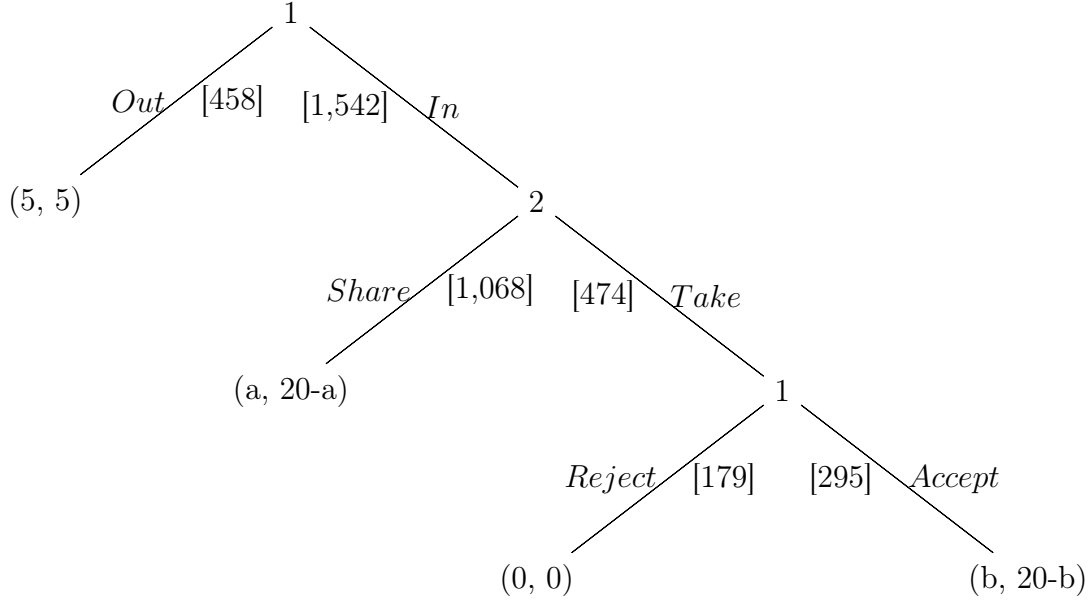


Figure 3. Summary of experimental results

differences reflect the differing incentives for Player 2 in the two types of games: the payoff to Player 2 from selecting *Grab* after *In* is 11 in game LT5 but only 5 in the High Take game.

We also observe a decreasing tendency for Player 1's to select *Reject* as b increases. When $b = 1$, the relative frequency of *Reject* choices is 75% for the Low Take game LT1 and 89% for the High Take game HT1. However, our standard errors are especially large in the LT1 game, where in only 4/200 instances did the game reach the third stage. The proportion of *Reject* choices was 14% in game LT5 and 17% in game HT5. In general, there appears to be a pattern of greater proportions of *Reject* choices in the High Take games, though again standard errors are large.

Beliefs and Plans. Figure 5 shows means of Player 1's self-reported plans and beliefs by game. In general, the patterns across two *Take* amounts and costs of *Reject* are similar to the observed choices. The data here is suggestive of bias: the mean self-reported *Out* plan is typically greater than the empirical probability of *Out* choices and the mean *Share* belief is mostly below the empirical probability of *Share* choices. In addition, the mean reported likelihood of choosing *Reject* is lower than the empirical probability of *Reject* in all but the games where $b = 5$. Next, we investigate the relation between beliefs and behavior more thoroughly.

Relating behavior, beliefs, and plans. Probability models are well-calibrated when probabilities match observed relative frequencies. For example, a weather forecasting model

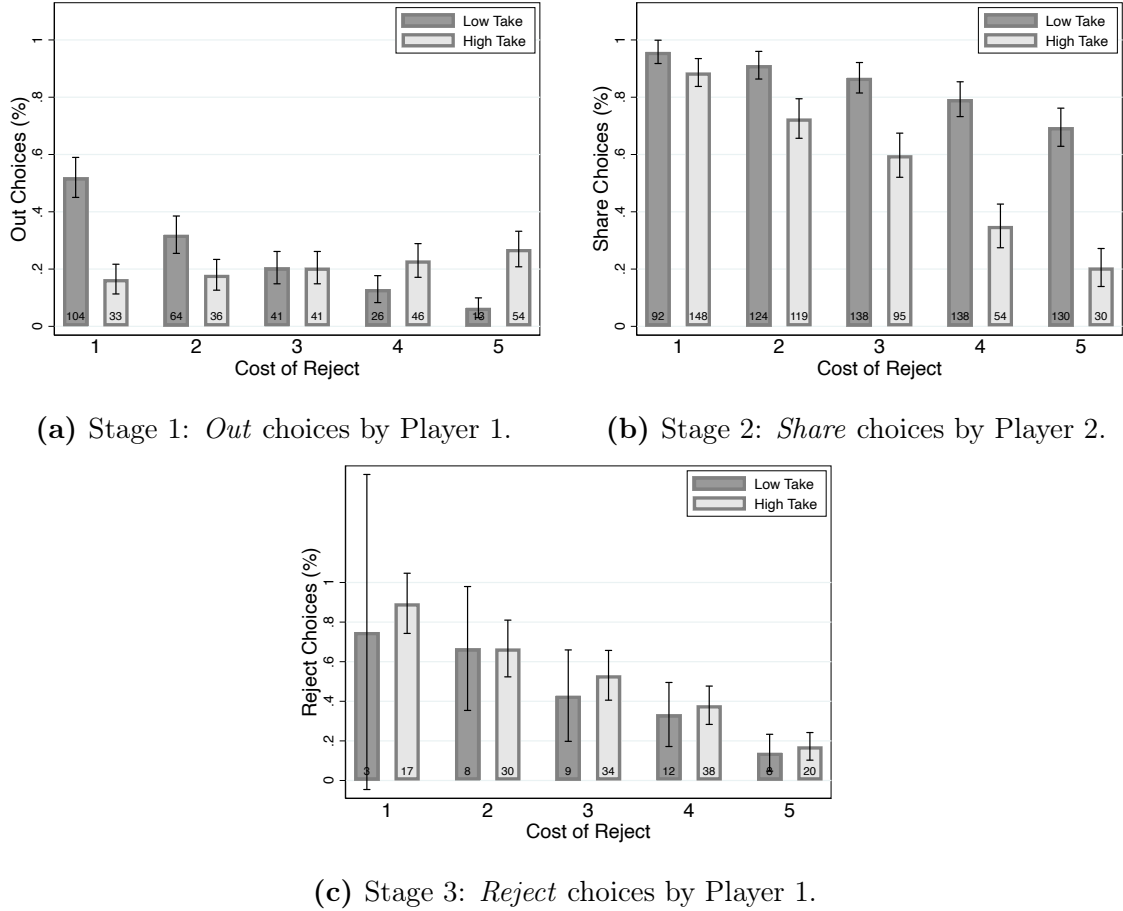


Figure 4. Relative frequency of player actions, by the cost of *Reject* and *Take* amount.

is well calibrated if, when the forecast indicates an 80% chance of rain, it rains 80% of the time. In addition probability models should also have good resolution, meaning that when the observed outcome is very likely (unlikely), the predictor variable is close to 1 (0). In the machine learning literature, it is commonplace to evaluate classifier performance by plotting Receiver Operating Characteristic curves. These curves plot the true positive rate (sensitivity) of a predictor versus the false positive rate (1 minus the specificity) as a function of the threshold that determines the prediction. One can then compute the Area Under the Curve (AUC) summary statistic, with values closer to 1 indicating better classification.¹⁶ Another approach is to estimate a simple linear probability model that regresses the observed outcome on the elicited probabilities. The fitted model will have slope 1 and intercept 0 if it is perfectly calibrated; the higher the R^2 of the regression, the greater the resolution of the prediction.

¹⁶See *e.g.* Murphy (2012, Chapter 5).

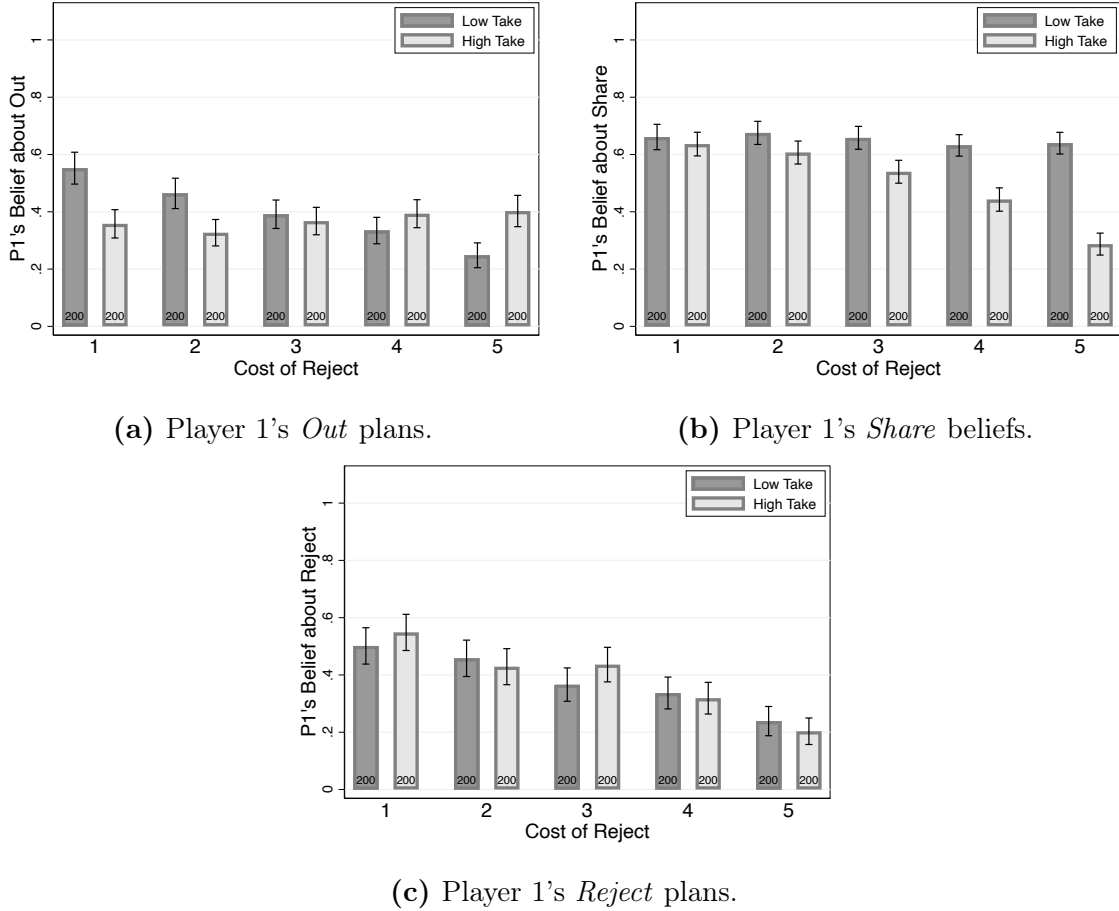


Figure 5. Player 1's beliefs and plans, by the cost of *Reject* and *Take* amount.

In the Appendix, we perform both types of analyses for each of Player 1's 3 self-reported probabilities of *Out*, *Share*, and *Reject* (Supplementary Figures 1, 2, and 3). Though clearly biased (Supplementary Figures 1(a) and 3(a)), elicited plans are good predictors of behavior, with R^2 values of 0.45 and 0.38 for *Out* and *Reject* plans (Supplementary Figures 1(a) and 2(a), respectively), and AUC statistics 0.9167 (*Out*, Supplementary Figure 1(b)) and 0.8370 (*Reject*, Supplementary Figure 3(b)). Elicited beliefs about co-player behavior are less accurate predictors, though still informative: a simple regression of *Share* choices by Player 2 on Player 1's beliefs has R^2 of 0.10, and the AUC statistic is 0.6828 (Supplementary Figures 2(a) and 2(b)).

4.2 Main Results

In this section we revisit our preconceived hypotheses. We start by evaluating the relationship between Player 1’s self-reported beliefs, plans, and choices (Hypothesis 1), followed by an examination of the relationship between game structure and behavior (Hypothesis 2). Then, we test for the effect of communication on behavior and payoffs (Hypothesis 3). Finally, we investigate how message content in the form of promises influences beliefs and behavior (Hypothesis 4 and 5).

4.2.1 Testing H1: The effect of beliefs about *Share* on *Reject* choices and plans

To study the role of beliefs in driving costly punishment, in Table 2 we show regression analyses that study the link between beliefs, *Reject* choices, and *Reject* plans. These analyses account for the cost of choosing *Reject*, the *Take* amount, and whether communication was available. In this subsection, we focus on the role of beliefs; in subsections 4.2.2 and 4.2.3 we discuss the effects of the *Take* amount and of communication.

Columns A-B in Table 2 report the results of logistic regressions where the dependent variable, P1’s *Reject* Choice, is equal to 1 if Player 1 *Rejects* the offer in stage 3, and equal to 0 if Player 1 *Accepts* the offer in stage 3. We find, in the model in Column B, a significant relationship between Player 1’s first order belief about *Share* and decision to *Reject* the offer after *Take*. A 10% increase in the elicited probability of *Share* increases Player 1’s chance of rejecting by 2.488%, which is consistent with Hypothesis 1.

The models in Columns A-B do not include subject or session level controls. When either of these controls are included, the coefficient on Belief about *Share* is not significant (see Supplementary Table 1), implying that the relation in Column B is driven by variability between individuals or sessions. This is not inconsistent with our theory, which simply posits a relation between beliefs and behavior. We also have limited data: the 474 *Take* choices imply that on average we observe 4.74 choices to *Accept* or *Reject* at the end of the game (min 1, max 10). Thus, our data may not be sufficient to establish a within-participant relation between Beliefs about *Share* and *Reject* choices.

In Table 2, Columns C-D, we employ fixed effects linear regressions to study the determinants of Player 1’s reported *Reject* plan (divided by 100, to scale between 0 and 1). Here, we have data for each game played, for a total of 2,000 observations (20 per participant in the

Table 2. The Effect of Belief about *Share* on P1’s *Reject* Choice and Plan.

	P1’s <i>Reject</i> Choice		P1’s <i>Reject</i> Plan	
	A	B	C	D
	mfX / se	mfX / se	mfX / se	mfX / se
Cost of <i>Reject</i>	-0.1985*** (0.0219)	-0.1814*** (0.0245)	-0.0736*** (0.0077)	-0.0677*** (0.0069)
High <i>Take</i>	0.0442 (0.0584)	0.0769 (0.0606)	0.0087 (0.0078)	0.0272*** (0.0082)
Communication	0.0174 (0.0494)	0.0111 (0.0483)	0.0552*** (0.0170)	0.0449*** (0.0167)
Period	0.0137*** (0.0044)	0.0129*** (0.0045)	0.0122*** (0.0014)	0.0118*** (0.0017)
Belief about <i>Share</i>		0.2488** (0.1173)		0.1230*** (0.0360)
Observations	474	474	2000	2000
BIC	560.4	558.9	589.4	571.6
Subject controls	No	No	Yes	Yes

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$. Standard errors in parentheses.

Note: mfx: marginal effect. Marginal effects for continuous variables are evaluated at means, and for binary variables are evaluated as the discrete change from 0 to 1. se: standard error. Standard errors are bootstrapped at the session level. Fixed effect logistic regressions are employed for P1’s *Reject* Choice, and fixed effect linear regressions are employed for P1’s *Reject* Plan. See Supplementary Table 1 and 2 for additional regressions of *Reject* Choice and *Reject* Plan.

role of Player 1), and so these models include participant level controls. Focusing on the role of beliefs, the model in Column D shows that Player 1’s first order Belief about *Share* has a positive and statistically significant association with Player 1’s *Reject* plan, consistent with the frustrated anger model and with Hypothesis 1. Thus, within-participant variation in the probabilistic belief that Player 2’s would choose *Share* was linked to changes in participants’ probabilistic plan to *Reject*: an increase of 10% in the Belief about *Share* is associated with about a 1.2% increase in the reported Plan to *Reject*.

4.2.2 Testing H2: The effect of the cost of *Reject* and the *Take* amount on choices and plans

Our model posits that participants make economic tradeoffs between monetary rewards frustrated anger, and that larger *Take* amounts should lead to increased frustration after *Grab*. Thus in our experimental design (see Table 1), we varied both the cost of *Reject* (b) and the *Take* amount ($a - b$, either 4 or 10).

The results in Figures 4(c) and 5(c) show declining rates of *Reject* choices and plans to *Reject* as b increases and are consistent with the idea that participants are sensitive to the cost (foregone reward) of choosing *Reject*. After pooling High and Low *Take* games, a nonparametric test for trend (Cuzick, 1985) rejects the null hypothesis that *Reject* rates are unrelated to the foregone payoff from *Accept* (b) ($p < 0.001$). A nonparametric test for trend also rejects the null hypothesis that *Reject* plans are unrelated to the cost of *Reject* ($p < 0.001$). Furthermore, the regression analyses in Table 2 show that the coefficient estimate on the variable “Cost of *Reject*” (b) is significant and negative for both *Reject* choices (Models A-B) and plans (Models C-D). The magnitude of the effect varies from an almost 20% increase per unit (Models A-B, choices) to about a 7% increase per unit (Models C-D, plans). Participants were clearly sensitive to the cost of choosing *Reject*, consistent with Hypothesis 2.

Our theory also implies a relationship between the *Take* amount ($a - b$) and *Reject* choices and plans, so next we return to the regression analyses in Table 2. In neither of the regression analyses for Player 1’s *Reject* choice is the coefficient on High *Take* significant (Models A-B), though the sign is positive. In the Appendix, Supplementary Table 1 reports the results from additional specifications; in only one of the models is the coefficient for High *Take* significant (Model D). Turning to our fixed-effects linear regression analyses for *Reject* plans (r_1), we again find that the coefficient on High *Take* is positive but not statistically significant in Model C. However, after including Player 1’s Belief about *Share*, Model D finds that the High *Take* condition adds about 2.7% to the reported likelihood that Player 1’s will assign to *Reject*. While consistent with Hypothesis 2, the effect of High *Take* is conditional on beliefs. Additional model specifications give similar results; see Supplementary Table 2.

We interpret these results through the lens of the preceding analyses indicating the importance of player’s beliefs about whether their co-players will cooperate. Clearly, Player 1s’ *Reject* choices and plans are closely linked to their beliefs about whether Player 2 will *Share*, and given beliefs, on the *Take* amount ($a - b$).

4.2.3 Testing H3: The effect of communication on cooperation

Each experimental session included both no-communication and communication blocks. In the latter, Player 2 was given the opportunity to send a pre-play free-form message to Player 1. We first investigate how the communication treatment affects game outcomes, and the result is shown in Table 3.

Table 3. The Effect of Communication.

	Out (<i>Out</i>)	Cooperation (<i>In, Share</i>)	Rejection ((<i>In, Reject</i>); <i>Take</i>)	Acceptance ((<i>In, Accept</i>); <i>Take</i>)	Total
No Communication	263 26.30%	467 46.70%	97 9.70% 35.93%	173 17.30% 64.07%	1000 100.00% 100.00%
Communication	195 19.50%	601 60.10%	82 8.20% 40.20%	122 12.20% 59.80%	1000 100.00% 100.00%
Total	458 22.90%	1068 53.40%	179 8.95% 37.76%	295 14.75% 62.24%	2000 100.00% 100.00%

Note: Row 1: number of observations; row 2: percent of total observations; row 3: percent of observations that reach the third stage.

The cooperative outcome (*In, Share*) is more prevalent in the communication treatment (60.10% vs. 46.70%). A 1-sided Fisher's exact test confirms that the cooperation rate is higher in the communication treatment (p-value < 0.001). This result is consistent with the belief dependent models of frustrated anger and guilt aversion and with Hypothesis 3, that communication will increase cooperation. A chi-squared test shows that allowing communication has a significant effect on the distribution of outcomes (terminal histories) (p-value < 0.001). The conditional *Reject* rate is also higher in the communication treatment (40.20% vs. 35.93%), but this difference is not significant (1-sided Fisher's exact test: p-value = 0.197).

Communication also affects reported beliefs. Figure 6 presents histograms of Player 1's self-reported plans for choosing *Out* and *Reject* and beliefs that player 2 will choose *Share*, in the communication and the no-communication treatments. In the communication treatment we measured beliefs both before and after messages were received; unless otherwise noted

the belief data we report for the communication treatment were recorded after messages were received. Epps-Singleton tests confirm that the distributions of reported beliefs are significantly different in the communication vs. the no-communication treatment (plan for *Out*: p-value < 0.001; belief for *Share*: p-value < 0.001; plan for *Reject*: p-value < 0.001).

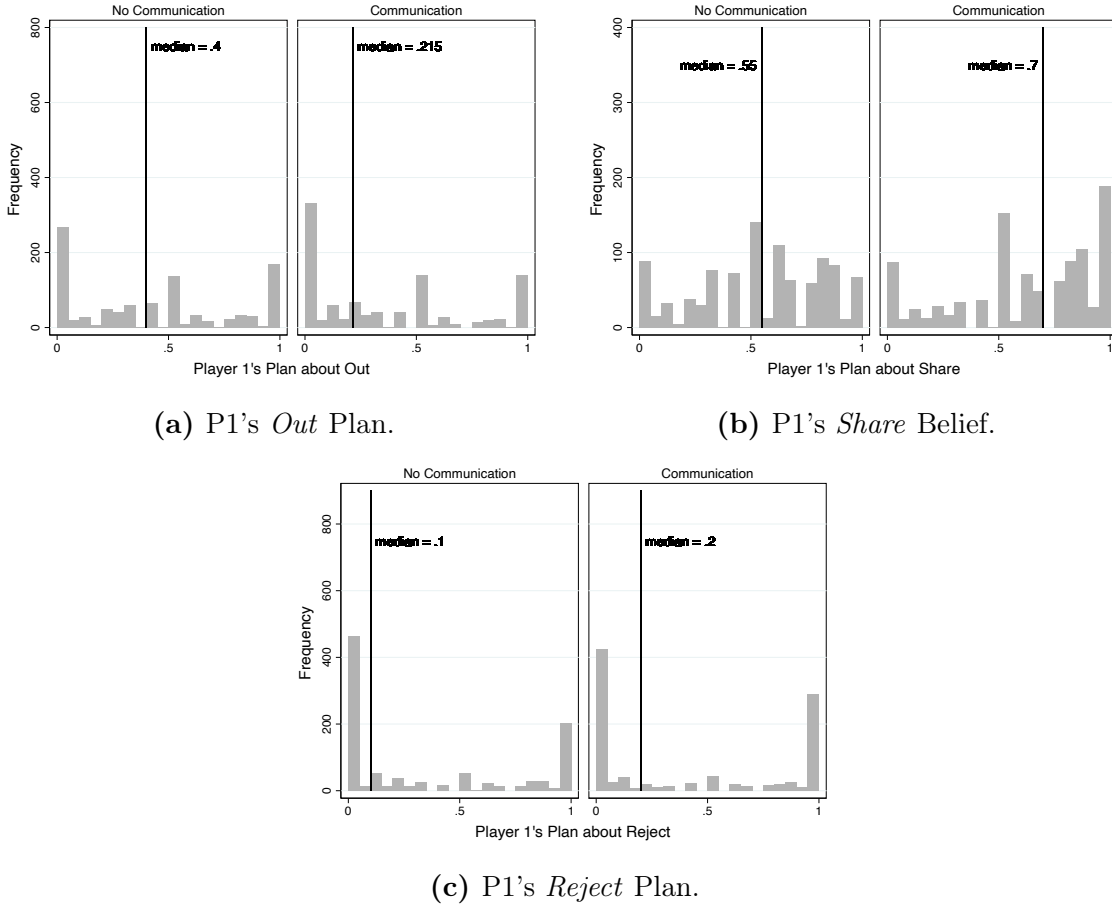
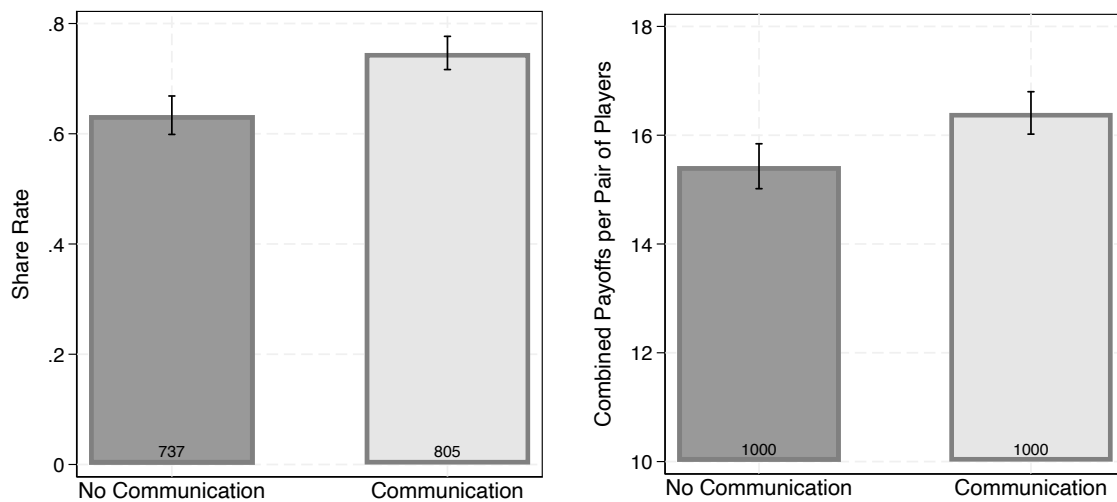


Figure 6. Histograms of P1's Plans (beliefs about own actions) and Beliefs (about P2), by Communication.

As predicted by Hypothesis 3, communication has a strong effect on efficiency and co-operation. Figure 7(a) compares *Share* outcomes from the no-message and message blocks, pooling the data from all sessions, with number of observations labeled on the bars. We observe a significantly higher cooperation rate with communication on a subject level (1-sided t-test: p-value < 0.001). Additionally, we show that in Supplementary Figure 4(a), relatively higher cooperation rates are observed across 10 different game variations in the communication treatment.

To test whether communication improves efficiency, we first compare participants' combined payoffs (Figure 7(b)). On average, combined payoffs are significantly higher in the



(a) P2's *Share Rate* with Communication.

(b) P1 and P2's Combined Payoffs.

Figure 7. Communication Improves Cooperation and Earnings.

communication treatment (\$16.41 vs. \$15.43, rank sum test: $p\text{-value} = 0.014$). Next, we look into Player 1 and Player 2's payoffs separately.

Figure 8 shows that on average, the payoffs of Player 2's are insignificantly greater in the communication treatment (\$9.03 vs. \$9.34, rank sum test: $p\text{-value} = 0.127$); whereas, Player 1's average earnings are significantly larger with communication (\$6.40 vs. \$7.07, rank sum test: $p\text{-value} < 0.001$). This suggests that social welfare or efficiency increases if communication is allowed. This result is consistent with our Hypothesis 3, that communication improves efficiency.

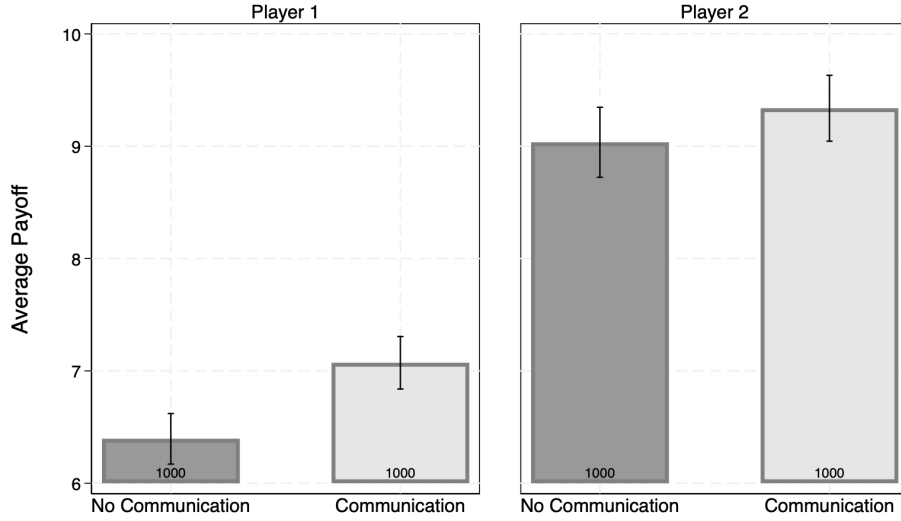


Figure 8. The Effect of Communication on Payoffs by Player Role.

4.2.4 H4 and H5: Communication Shifts Beliefs through Promises

We now consider the effect of communication on beliefs (measured after the message was received in the communication treatment). Figure 9 shows that communication affects Player 1's reported beliefs, and this result is consistent with our Hypothesis 4. Player 1's report a higher likelihood that Player 2 will cooperate (1st order belief about *Share*, q_1) when communication is allowed. Communication affects Player 1's own plans as well. With communication, Player 1 believes that she is less likely to play *Out* but more likely to *Reject* if 3rd stage is reached. 1-sided t-tests confirm that Player 1's beliefs are significantly different in the communication treatment and the no-communication treatment (plan for *Out*: p-value = 0.009; 1st order belief about *Share*: p-value < 0.001; plan for *Reject*: p-value = 0.069). In addition, the direction of how communication influences beliefs is consistent with belief-dependent anger.

To examine the links between message content, beliefs, and behavior, we manually coded messages as promises if they follow the pattern of "If you choose *In*, I will choose *Share*."¹⁷ Using this approach, we identify 32% of messages as promises, and the median number of promises per session was 32.2%. Sample messages and their categories are presented in Supplementary Table 6.

As noted above, in the communication treatment, we measured beliefs both before and

¹⁷As a reminder, we used neutral action labels in the actual experiment.

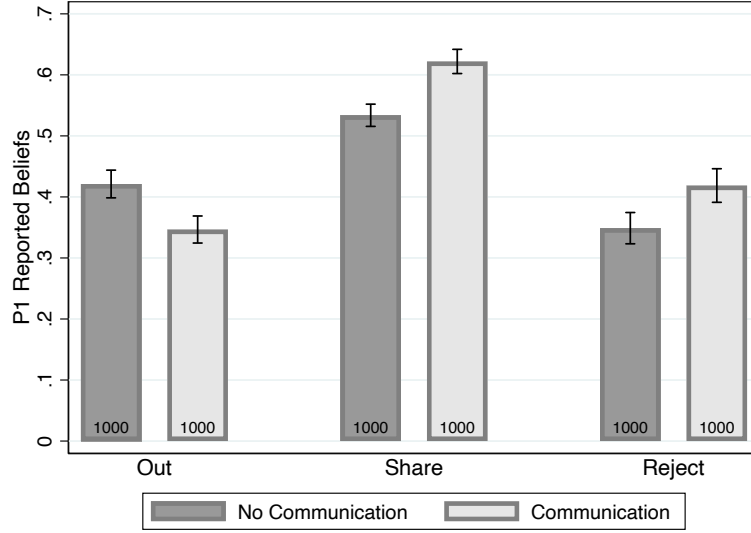


Figure 9. Communication Influences P1's Reported Beliefs.

after receiving a message. Figure 10 shows that promises have a strong effect on Player 1's reported beliefs. Promises increase Player 1's belief about Player 2's cooperative behavior (1st order belief about *Share*). Promises also influence Player 1's beliefs about their own actions (plan for *Out* and *Reject*). After receiving a promises message, Player 1s report that they will be less likely to choose *Out*, but will be more likely to punish Player 2. 1-sided t-tests show a significant difference in the change in Player 1's reported beliefs after receiving a promise, compared to receiving a message that did not involve a promise (plan for *Out*: p-value < 0.001; 1st order belief about *Share*: p-value < 0.001; plan for *Reject*: p-value = 0.011). In addition, two-sided t-tests confirm that the change in Player 1's reported beliefs after non-promise messages is not significantly different from 0 (plan for *Out*: p-value = 0.779; 1st order belief about *Share*: p-value = 0.343; plan for *Reject*: p-value = 0.390). Promises have a significant effect upon beliefs, while non-promises have an insignificant effect, consistent with Hypothesis 4.

4.2.5 Promises Influence Behavior

To further demonstrate the effect of promises on behavior as predicted in Hypothesis 5, we look at behavior differences under promises and non-promises. Supplementary Table 3 shows the outcome distribution with respect to promises and non-promises is consistent with belief-dependent anger. A chi-squared test shows that the distribution of outcomes is significantly different with and without promises (p-value < 0.001). Figure 11 shows that the

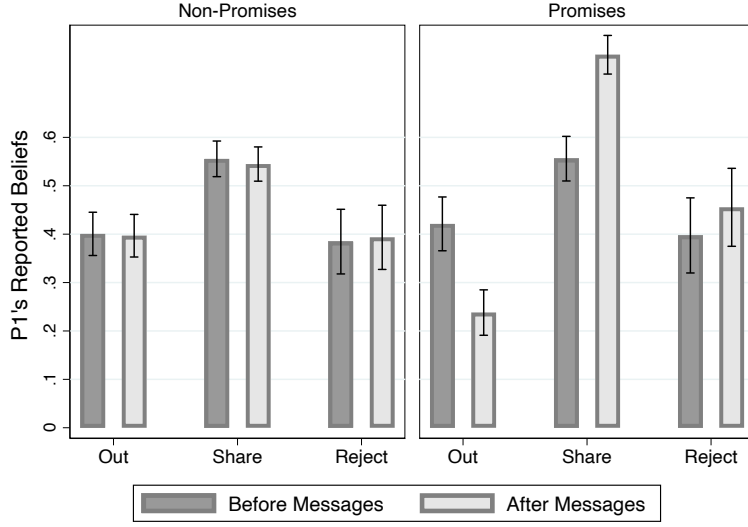
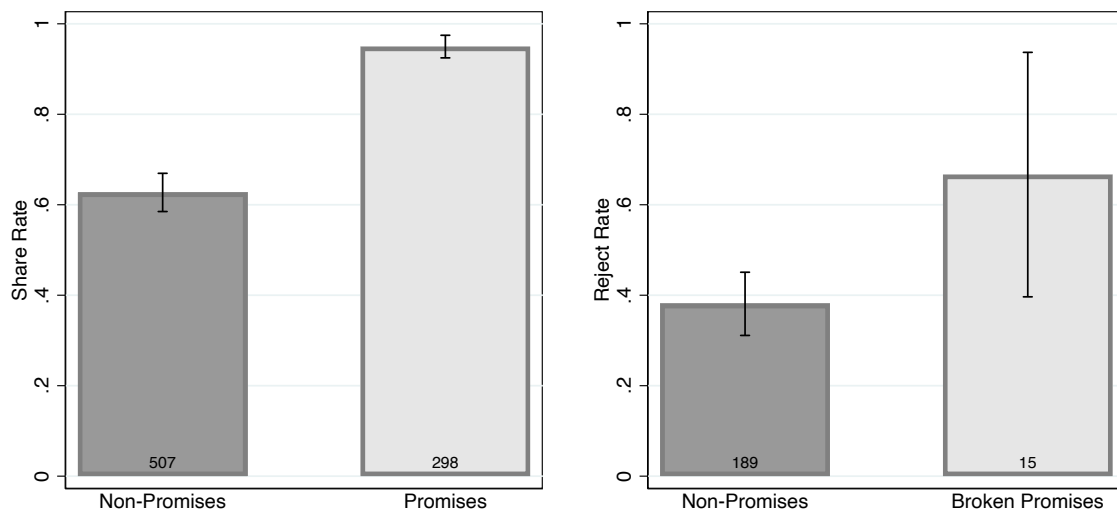


Figure 10. Belief Change After Receiving a Message.

proportion of both *Share* and *Reject* choices is higher when a promise is made. This result is consistent with Hypothesis 5: promises foster cooperation, but broken promises lead to higher rates of punishment. The effect of promises is greater than the effect of communication, see the difference between Supplementary Figure 4(a) and 4(b), that a greater improvement of cooperation can be observed across 10 game variations when separating promise and non promise messages than with communication treatment alone.

The result shown in Figure 11(a) is consistent with the frustration-anger model, in that if Player 2 anticipates the change in Player 1's beliefs following a promise, Player 2 will have increased motivation to choose *Share* in order to avoid punishment from a *Reject* choice after *Take*. When we compare Player 2's behavior after non-promises vs. after promises, the *Share* rate is significantly higher following promises (1-sided Fisher's exact test: $p\text{-value} < 0.001$). A rank sum test confirms that subject level *Share* rate is also higher with promises ($p\text{-value} < 0.001$). This result holds on the level of each of the 10 different game variations as well. Supplementary Figure 4(b) shows that the *Share* rate for promises is higher across all 10 game variations. Promises affect not only cooperative behavior but also rejection. As predicted by the frustration-anger model, Player 1's beliefs change following promises, and Player 1 is more likely to punish with broken promises. Figure 11(b) shows that the *Reject* rate is higher with broken promises compared to non-promises, consistent with Hypothesis 5 (1-sided Fisher's exact test: $p\text{-value} = 0.030$). A rank sum test confirms that subject level *Reject* rate is also higher with broken promises ($p\text{-value} = 0.068$).



(a) P2's *Share* Rate with Promises.

(b) P1's *Reject* Rate with Broken Promises.

Figure 11. Kept and Broken Promises.

4.3 Additional Observations

This section reports additional observations regarding our data. We first discuss the observed persistent effect of communication, then we present evidence for gender differences.

4.3.1 The Persistent Effect of Communication

There is a significant difference between the experimental sessions with communication first and the sessions with communication second. Figure 12 shows that in the communication-first sessions, there is a persistent effect of communication: the higher rate of cooperation (*Share* outcomes) is sustained in the subsequent no-communication periods. Restricting attention to the first 10 rounds of each treatment, there is a significantly higher cooperation rate in the communication-first sessions than in the communication-second treatment (58.86% vs. 35.18%, 1-sided Fisher's exact test: $p\text{-value} < 0.001$). The difference disappears in rounds 11-20 (61.36% vs. 61.07%, 2-sided Fisher's exact test: $p\text{-value} = 0.948$). This suggests that the communication effect is so strong that after being exposed to the communication environment, participants behave as if they are still sending and receiving messages in the second, no-communication, block. This durable effect of communication arises not just in outcomes, but also in Player 1's reported beliefs. Player 1's have a significantly higher first order belief about *Share* in the first 10 periods of the communication-first sessions relative to

the no-communication sessions, (two-sided t-test for difference in session means: $t = 3.628$, $df = 9$, p-value = 0.006), but there is no difference in beliefs about *Share* when comparing rounds 11-20 of the communication-first and communication-second sessions ($t = -0.416$, $df = 9$, p-value = 0.6874).

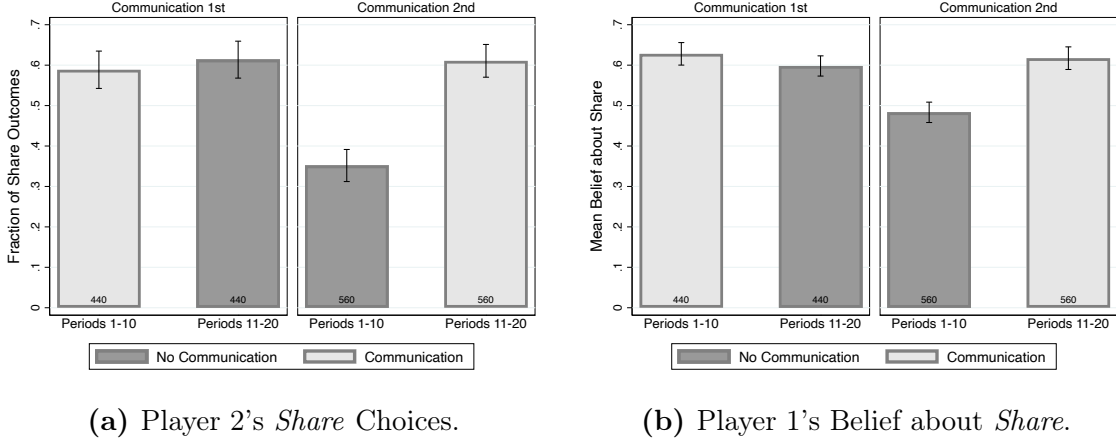


Figure 12. Persistent Communication Effect.

Because of this persistent effect of communication, we examine the distribution of outcomes after restricting the sample to include only the first 10 rounds. Figure 13 compares the effects of communication on the distribution of outcomes with all 20 rounds and with the first 10 rounds only, when the no-communication group has no experience with messages. The contrast of communication vs. no communication is stronger when we look at the first 10 rounds only. The mean fraction of *Share* outcomes in the communication treatment in the first 10 rounds is 58.86%, which is close to the overall mean for 20 rounds (60.10%, see also in Table 3), but the cooperation rate without communication in the first 10 rounds decreases to 35.18%. A chi-squared test shows that the communication treatment has a significant effect on the distribution of outcomes for the first 10 rounds of the experiment (p-value < 0.001). These results demonstrate that communication has a strongly positive and persistent effect on cooperation.

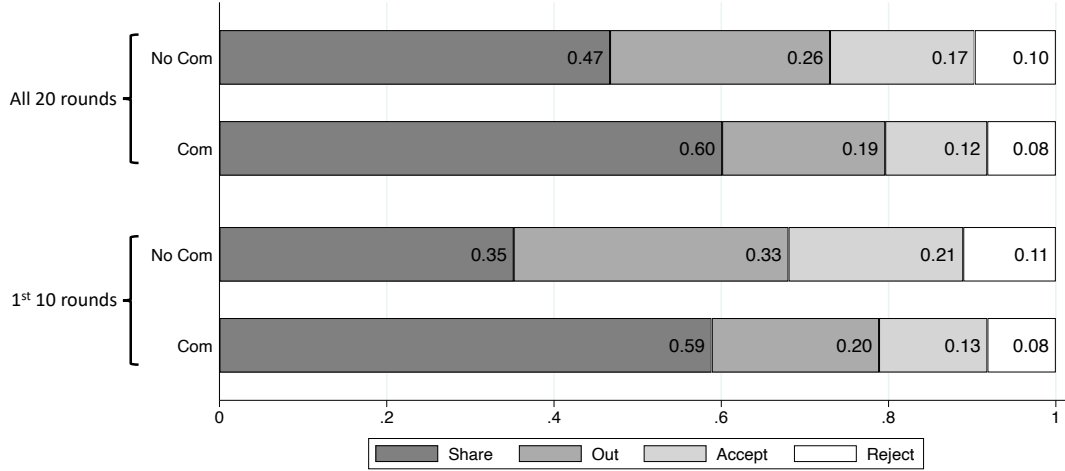


Figure 13. Outcomes Compare All 20 vs. 1st 10 Rounds.

4.3.2 Gender Differences

We started this project with BDS’ theory in mind and the intention to test hypotheses 1-5. We recorded subjects’ genders without any preconceived conjectures as regards whether results would differ between women and men. Aina et al. (2020), however, report that men are more affected by anger than women. As it turns out, we have comparable findings. Females and males’ behavior are relatively similar, except that when promises are made, males tend to *Reject* more often (*Reject* rate 70% vs. 30%, rank sum test: p-value = 0.077). We ran the fixed effects linear regressions for *Reject* plans separately for females and males and report the results in Table 4. The effect of communication survives with females, but disappears with males. Whereas, beliefs about *Share* are significantly positively associated with plans to *Reject* for males, but not for females. Consistent with Aina et al. (2020), the coefficient estimates for “High *Take*” and “Belief about *Share*” are positive and statistically significant predictors of Player 1’s *Reject* plan in the regression for males, but not in the female-only analysis. In our data as well, men’s beliefs and choices are more consistent with the frustration-anger model.

Table 4. Linear Regressions – Gender Effect of P1’s *Reject* Plan.

	Females		Males	
	A	B	C	D
	mfx / se	mfx / se	mfx / se	mfx / se
Cost of Punishment	-0.0744*** (0.0144)	-0.0703*** (0.0135)	-0.0732*** (0.0083)	-0.0668*** (0.0087)
High <i>Take</i>	-0.0027 (0.0138)	0.0104 (0.0148)	0.0165 (0.0161)	0.0365** (0.0157)
Period	0.0110*** (0.0029)	0.0107*** (0.0031)	0.0130*** (0.0030)	0.0126*** (0.0030)
Communication	0.0592* (0.0333)	0.0539 (0.0351)	0.0543* (0.0319)	0.0397 (0.0318)
Belief about <i>Share</i>		0.0734 (0.0714)		0.1539*** (0.0308)
Observations	880	880	1100	1100
BIC	154.8	157.5	451.8	437.6
Subject controls	Yes	Yes	Yes	Yes

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$. Standard errors in parentheses.

Note: mfx: marginal effect. Marginal effects for continuous variables are evaluated at means, and for binary variables are evaluated as the discrete change from 0 to 1. se: standard error. Standard errors are bootstrapped at the session level. Fixed effect linear regressions are employed for P1’s *Reject* Plan. Supplementary Table 4 and 5 for additional regressions of *Reject* Plan separating for females and males.

5 Alternative Theories of Motivation

Experimental and behavioral economists have convincingly argued that models of social preferences are needed to explain human behavior, but little of such work factors in anger and frustration. One may wonder if doing so is necessary. For example, can models of inequity aversion (e.g. Fehr and Schmidt, 1999; Bolton and Ockenfels, 2000) explain our data? One implication of inequity aversion is that if player 1 ever *Rejects* a high offer in the 3rd stage, then she/he would never *Accept* a lower offer, regardless of communication or beliefs. Using this idea we can classify subjects into four categories, shown in Table 5. “IA Violation” represents subjects whose behavior is inconsistent with inequity aversion: they either *Reject* a higher and *Accept* a lower offer, or they both *Reject* and *Accept* the same offer (e.g. rejecting an offer of 3 in one round and accepting 3 in another). “Inequity

Averse” subjects’ behavior is always consistent with inequity aversion, “Self-interest” refers to players who always *Accept* any offer, and “Unclassified” are subjects that faced fewer than two different offers.

Table 5. Classification of Player 1 behavior.

	IA violation	Inequality averse (IA)	Self-interest	Unclassified
# of Subjects	36	28	33	3
# of 3rd Stage Decisions	5.42	4.79	4.27	1.33

Table 5 indicates that 36% of subjects are inconsistent with either self-interest or inequity aversion, 28% of subjects demonstrate behavior consistent with inequity aversion, while 33% of subjects behave as if they care only for material self-interest. Moreover, the number of subjects whose behavior is inconsistent with inequity aversion or self-interest increases when subjects have more decisions in the 3rd stage. This suggest that inequity aversion cannot explain the behavior of at least one-third of our participants, and models that allow for non-consequential behavior such as BDS may be needed to fully capture the range of behavior demonstrated.

Another strand of models addresses subjects’ tendency to honor promises, whether they be motivated by belief-dependent guilt aversion (see Charness and Dufwenberg (2006) and Battigalli and Dufwenberg (2007)) or the direct preference to honor a promise (e.g. Vanberg, 2008). These approaches help explain why communication increases the frequency of *Share* choices, but our results indicate that (the avoidance of) frustration, anger, and costly punishment in our game has additional effects. First, models of a tendency to honor promises and belief-dependent guilt aversion cannot explain the behavioral results we observe in the third stage of the games, regarding increased rates of punishment when promises are breached. Second, after promises, participants in our study choose to *Share* a striking 95% of the time. This amount of promise-keeping is much higher than in comparable studies without a punishment stage (e.g. Charness and Dufwenberg, 2006).

An alternative motivation for punishing broken promises might involve reciprocity, in which individuals are motivated to reward kindness and punish hostility through beliefs. However, sequential reciprocity (Dufwenberg and Kirchsteiger, 2004) allows for the possibility that players engage in mutual unkindness on the equilibrium path. In our setting, this would imply that a Player 1 whose assigned a low probability to Player 2 choosing *Share* could also

report plans (beliefs about her own subsequent actions) that involve a high probability of choosing *Out*, and a low probability of choosing *Reject*. This pattern of beliefs is ruled out by the frustration and anger model, and in fact, we do not observe it in our data. Player 1s' reported beliefs about the likelihood of *Share* are negatively associated with self-reported plans to choose *Out*, and positively associated with plans to choose *Reject* (Figure 14).¹⁸ This observation does not constitute a refutation of reciprocity theory, as it may be possible to observe such mutual unkindness in other settings (e.g. feuds), but it does suggest that the beliefs and plans we elicit are consistent with the frustration-anger model.

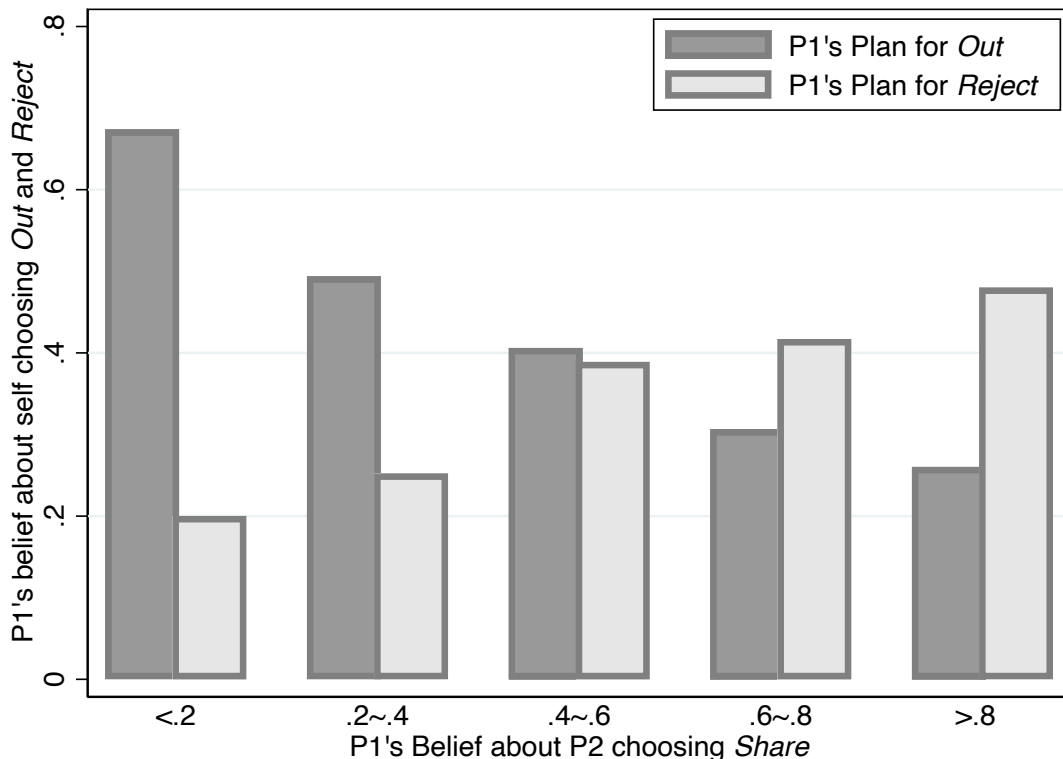


Figure 14. Player 1's self reported plans vs. beliefs about *Share*.

Another potential motivation is the desire to conform to social norms. Many individuals are willing to pay a cost to punish norm violators (Fehr and Gächter, 2000; Xiao and Houser, 2011). In our setting, Player 1 might regard Player 2's decision to *Take* as a violation of fairness norms. Player 1 can then choose to punish this norm violation. The frustration-anger explanation and the punishing-norm-violators are not mutually exclusive, as frustration is likely to result from the violation of the fairness norm. Krupka et al. (2017) show that social (in)appropriateness ratings (norms) and second-order beliefs (guilt aversion) are both good

¹⁸Bosman and Van Winden (2002) observed a similar pattern in the context of a power-to-take game.

predictors of behavior involving informal agreements. This suggests that further work is necessary to separately identify the motivations of belief-dependent anger from the desire to punish norm violators in the context of broken promises.

6 Conclusion

We study the effect of communication on strategic behavior in environments that allow for trust, promises, deception, and punishment. Communication increases cooperation and impacts beliefs. Beliefs are shaped by promises, and punishment increases with broken promises. The results support the idea that communication, beliefs, and costly punishment are linked through the mechanism of belief-dependent frustration and anger.

The hold-up problem arises from relationship-specific investments and the lack of verifiability of contracts. It is a central challenge in the economics of contracting and organizations. While other studies of hold-up focus on incomplete information and aspects of the bargaining process, our simple model emphasizes that hold-up problems arise when it is not possible to fully commit to uphold a contract. We show how the hold-up problem can be resolved at least in part by modeling emotional agents who are prone to anger.

Communication from the second-mover to the first mover may be beneficial in our setting for two reasons. First, there are two sequential equilibria of the hold-up game with angry second movers, and communication may help equilibrium selection. Second, communication may influence players' beliefs. We show that promises have a strong impact on player's belief that their co-players will share. An increase in the probability of sharing raises players' expected payoff, resulting in greater frustration and an increased propensity to engage in costly punishment when promises are broken.

Our analysis emphasizes that anger and the threat of costly punishment can help with the hold-up problem, and that communication can further aid in facilitating cooperation. In limiting our analysis to the emotion of anger we rule out many related concerns, including inequality aversion, guilt aversion, reciprocity, and concern for social norms. All these factors may play a role in informal contracting. Future work will shed more light on when and how much each factor contributes to cooperation.

One way to think about our work is that we formalize ideas from Selten (1978), Hirshleifer (1987), Frank (1988), and others who suggest that emotions solve commitment problems.

Our belief-dependent agents are frustrated when their payoffs do not meet expectations. Building on the frustration-aggression hypothesis from psychology, we link frustration to anger and aggression. This model also captures the notion from evolutionary psychology that the function of anger is to resolve bargaining conflicts in favor of the angry individual (Sell et al., 2009).

Our solution to the hold-up problem has some similarities with the approach of Hart and Moore (2008, H&M). They model agents for whom contracts serve as a reference point by which outcomes are evaluated. Parties who feel shortchanged relative to the reference contract feel “aggrieved” and shade on performance. Our focus is different, in that the starting point for H&M is the contract as a reference point, ours is the belief-dependent model of frustration and anger of BDS. In the BDS setting, the “reference point” against which outcomes are judged is expected payoffs. We study how anger and communication can help to resolve the hold-up problem, while H&M focus on how behavioral considerations generate tradeoffs between flexible and rigid contracts.

In studying the role of communication, we have limited our attention to messages from the second-mover to the first-mover. This is restrictive. In particular, messages from the first-mover to the second-mover could involve threats that gain credibility with anger and frustration in the picture. The topic is so interesting that it warrants its own research exercise, which, in fact, we run as a companion project (Dufwenberg, Li, and Smith, 2025).

Appendices

A Instructions

Below is an example of the instructions for sessions with the communication treatment before the no communication treatment. The instructions for the second part of the experiment were given to all the subjects after the communication block was completed.

Part I Instructions

Welcome to the experiment. The purpose is to study how people make decisions in a particular situation. Please do not speak to other participants during the experiment. Feel free to ask a question at any time by raising your hand.

You will receive \$5 for participating. You have the potential to earn additional money based on your own and others' decisions, as described below. Your decisions and payoffs will remain confidential. You will be paid individually and privately, in cash, at the end of the experiment.

There are two parts to the experiment. Both parts consist of multiple rounds of simple games that will be described below. The order in which choices are made in the games will remain the same in each round, but the payoff to different actions may change, so please pay careful attention to the payoffs in each round. At the end of the experiment, you will be privately paid for one randomly selected round from the entire experiment.

At the beginning of the experiment you will be randomly assigned to the role of either Player 1 or Player 2, and your role will not change throughout the experiment. In each round you will be randomly matched with another person in the room to play the game.

Prior to the start of each round, Player 2 will have the option to send messages to Player 1 (maximum 140 characters). Player 2 may say anything that he or she wishes in this messages, with one exception: no one is allowed to identify him or herself by name or number or gender or appearance. Violations of this rule may result in the loss of Player 2's payment for that part of the experiment (experimenter discretion). In that case the paired Player 1 will receive the average amount received by other Player 1's in this session.

Please raise your hand now if you have any questions. Select Continue when you are ready.

The game consists of three stages. The picture below may help and will be shown in each round. Payoffs will change in each round, so please familiarize yourself with the picture. Player 1's payoffs are listed above Player 2's payoffs. The game proceeds as follows:

- Player 1 goes first and must decide between A and B.
 - If A is chosen, the game ends and both players receive \$5.
 - If B is chosen, the game proceeds to stage 2.
- If Player 1 chooses B, Player 2 must decide between C and D.
 - If C is chosen, the game ends with payoffs specified for that round.
 - If D is chosen, Player 1 will make another decision.

- If Player 2 chooses D, Player 1 will decide between E and F.
 - If E is chosen, the game ends and both players receive \$0.
 - If F is chosen, the game ends with payoffs specified for that round.

Please raise your hand now if you have any questions. Select Continue when you are ready.

In each game you will be asked to guess how likely it is that certain events (decisions made by you or the other player) will happen. Your response is very important to our research. You will be asked to state the percent chance that each event will happen. You may select any number between 0 and 100, with the number you select indicating the likelihood of the event occurring (100 = certain the event will happen, 0 = certain the event will not happen). You will be rewarded with \$5 for answering these questions. You have an option to choose to pledge to answer the guessing questions to the best of your knowledge by checking the box below:

☐ **By checking this box, I pledge that I will answer all guessing questions to the best of my knowledge.**

Please raise your hand now if you have any questions. Select Continue when you are ready.

Part II Instructions

Thank you for completing the first part of the experiment. In the second part of the experiment your assigned role will not change. The second part of the experiment is like the first part, with one change: no messages will be exchanged. As before, this part consists of multiple rounds. In each round you will be randomly matched with another person in the room to play the game.

The only difference from the first part is that no messages will be exchanged for the second part of the experiment.

Please raise your hand now if you have any questions. Select Continue when you are ready.

As before, the game consists of three stages. The picture below may help and will be shown in each round. Payoffs will change in each round, so please familiarize yourself with

the picture. Player 1's payoffs are listed above Player2's payoff. The game proceeds as follows:

- Player 1 goes first and must decide between A and B.
 - If A is chosen, the game ends and both players receive \$5.
 - If B is chosen, the game proceeds to stage 2.
- If Player 1 chooses B, Player 2 must decide between C and D.
 - If C is chosen, the game ends with payoffs specified for that round.
 - If D is chosen, Player 1 will make another decision.
- If Player 2 chooses D, Player 1 will decide between E and F.
 - If E is chosen, the game ends and both players receive \$0.
 - If F is chosen, the game ends with payoffs specified for that round.

Please raise your hand now if you have any questions. Select Continue when you are ready.

B Supplementary Tables and Figures

Supplementary Table 1. Determinants of P1's *Reject* Choice (Robustness Checks).

	A	B	C	D	E
	mfX / se	mfX / se	mfX / se	mfX / se	mfX / se
Payoff from <i>Accept</i>	-0.1985*** (0.0219)	-0.1814*** (0.0245)	-0.0806** (0.0411)	-0.1680*** (0.0283)	-0.1604* (0.0915)
High Take	0.0442 (0.0584)	0.0769 (0.0606)	0.0170 (0.0349)	0.1407** (0.0717)	0.0728 (0.0978)
Communication	0.0174 (0.0494)	0.0111 (0.0483)	0.0053 (0.0234)		
Period	0.0137*** (0.0044)	0.0129*** (0.0045)	0.0058 (0.0040)	0.0116 (0.0071)	0.0155 (0.0265)
Belief about <i>Share</i>		0.2488** (0.1173)	0.0684 (0.0934)	0.4949*** (0.1397)	0.3771 (0.3470)
Promise				0.0692 (0.2282)	0.0564 (0.5434)
Observations	474	474	474	204	204
AIC	539.6	533.9	462.9	230.2	186.1
BIC	560.4	558.9	483.7	250.1	202.7
Session controls	No	No	Yes	No	Yes
Subject controls	No	No	No	No	No

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$. Standard errors in parentheses.

Note: Logistic regressions with P1's *Reject* choice as the dependent variable. mfx: marginal effect. Marginal effects for continuous variables are evaluated at means, and for binary variables are evaluated as the discrete change from 0 to 1. se: standard errors. Standard errors are clustered at the session level.

Supplementary Table 2. Determinants of P1's *Reject* Plan (Robustness Checks).

	A	B	C	D	E
	coef / se	coef / se	coef / se	coef / se	coef / se
Payoff from <i>Accept</i>	-0.0723*** (0.0082)	-0.0736*** (0.0082)	-0.0677*** (0.0068)	-0.0741*** (0.0078)	-0.0690*** (0.0086)
High Take	0.0068 (0.0109)	0.0087 (0.0080)	0.0272*** (0.0080)	0.0208 (0.0158)	0.0377*** (0.0105)
Communication	0.0698* (0.0362)	0.0552*** (0.0171)	0.0449** (0.0198)		
Period		0.0122*** (0.0015)	0.0118*** (0.0017)	0.0196*** (0.0039)	0.0196*** (0.0038)
Belief about <i>Share</i>			0.1230*** (0.0361)	0.1122** (0.0530)	0.2232*** (0.0527)
Promise				0.0064 (0.0219)	-0.0089 (0.0325)
Constant	0.5623*** (0.0557)	0.4445*** (0.0505)	0.3561*** (0.0433)	0.3408*** (0.0834)	0.2527*** (0.0830)
Observations	2000	2000	2000	1000	1000
AIC	682.198	561.399	538.041	121.620	1062.800
BIC	704.601	589.403	571.646	151.067	1092.247
Session controls	No	No	No	No	Yes
Subject controls	Yes	Yes	Yes	Yes	No

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$. Standard errors in parentheses.

Note: Linear regressions with *Reject* plan as dependent variable. mfx: marginal effect. Marginal effects for continuous variables are evaluated at means, and for binary variables are evaluated as the discrete change from 0 to 1. se: standard errors. Standard errors are clustered at the session level.

Supplementary Table 3. The Effect of Promises on Outcomes (Communication Treatment Only).

	Out (<i>Out</i>)	Cooperation (<i>In, Share</i>)	Rejection (<i>(In, Reject); Take</i>)	Acceptance (<i>(In, Accept); Take</i>)	Total
Promises	22 6.88%	283 88.44%	10 3.12% 66.67%	5 1.56% 33.33%	320 100.00% 100.00%
Non-Promises	173 25.44%	318 46.76%	72 10.59% 38.10%	117 17.21% 61.90%	680 100.00% 100.00%
Total	195 19.50%	601 60.10%	82 8.20% 40.20%	122 12.20% 59.80%	1000 100.00% 100.00%

Note: Row 1: number of observations; row 2: percent of total observations; row 3: percent of observations that reach the third stage.

Supplementary Table 4. Determinants of P1's *Reject* Plan (Females)

	A	B	C	D	E	F
	coef / se	coef / se	coef / se	coef / se	coef / se	coef / se
Payoff from <i>Accept</i>	-0.0726*** (0.0147)	-0.0744*** (0.0144)	-0.0703*** (0.0132)	-0.0751*** (0.0160)	-0.0717*** (0.0171)	-0.0717*** (0.0174)
High Take	-0.0027 (0.0129)	-0.0027 (0.0125)	0.0104 (0.0159)	0.0006 (0.0192)	0.0099 (0.0240)	0.0106 (0.0216)
Communication	0.0742* (0.0416)	0.0592* (0.0323)	0.0539 (0.0373)			
Period		0.0110*** (0.0029)	0.0107*** (0.0032)	0.0253*** (0.0040)	0.0253*** (0.0041)	0.0253*** (0.0037)
Belief about <i>Share</i>			0.0734 (0.0711)		0.0609 (0.0927)	0.0668 (0.0949)
Promise				0.0241 (0.0255)	0.0098 (0.0179)	0.0031 (0.0251)
Constant	0.5172*** (0.0914)	0.4143*** (0.0790)	0.3572*** (0.0705)	0.3058*** (0.0939)	0.2575*** (0.0958)	0.2557*** (0.1115)
Observations	880	880	880	440	440	440
AIC	179.545	130.935	128.772	13.019	13.474	438.206
BIC	198.665	154.835	157.451	33.453	37.995	462.727
Session controls	No	No	No	No	No	Yes
Subject controls	Yes	Yes	Yes	Yes	Yes	No

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$. Standard errors in parentheses.

Note: Linear regressions with *Reject* plan as dependent variable for females subjects only. mfx: marginal effect. Marginal effects for continuous variables are evaluated at means, and for binary variables are evaluated as the discrete change from 0 to 1. se: standard errors. Standard errors are clustered at the session level.

Supplementary Table 5. Determinants of P1's *Reject* Plan (Males)

	A	B	C	D	E	F
	coef / se	coef / se	coef / se	coef / se	coef / se	coef / se
Payoff from <i>Accept</i>	-0.0724*** (0.0103)	-0.0732*** (0.0090)	-0.0668*** (0.0082)	-0.0823*** (0.0085)	-0.0773*** (0.0086)	-0.0740*** (0.0077)
High Take	0.0128 (0.0192)	0.0165 (0.0157)	0.0365** (0.0152)	0.0065 (0.0204)	0.0289 (0.0237)	0.0428* (0.0238)
Communication	0.0708 (0.0498)	0.0543 (0.0341)	0.0397 (0.0319)			
Period		0.0130*** (0.0032)	0.0126*** (0.0029)	0.0148*** (0.0056)	0.0150*** (0.0055)	0.0150*** (0.0050)
Belief about <i>Share</i>			0.1539*** (0.0298)		0.1527*** (0.0490)	0.2427*** (0.0481)
Promise				0.0335 (0.0354)	0.0069 (0.0338)	-0.0010 (0.0360)
Constant	0.6018*** (0.0743)	0.4742*** (0.0742)	0.3703*** (0.0582)	0.5308*** (0.1082)	0.4174*** (0.1102)	0.3476*** (0.0890)
Observations	1100	1100	1100	550	550	550
AIC	493.885	426.762	407.562	126.671	117.458	524.642
BIC	513.897	451.777	437.580	148.220	143.318	550.501
Session controls	No	No	No	No	No	Yes
Subject controls	Yes	Yes	Yes	Yes	Yes	No

* $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$. Standard errors in parentheses.

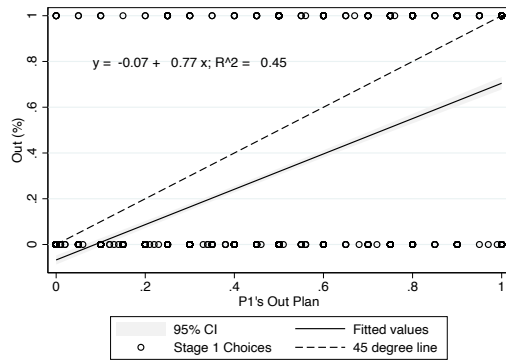
Note: Linear regressions with *Reject* plan as dependent variable for male subjects only. mfx: marginal effect. Marginal effects for continuous variables are evaluated at means, and for binary variables are evaluated as the discrete change from 0 to 1. se: standard errors. Standard errors are clustered at the session level.

Supplementary Table 6. Sample Messages. The full list of messages will be made available after publication.

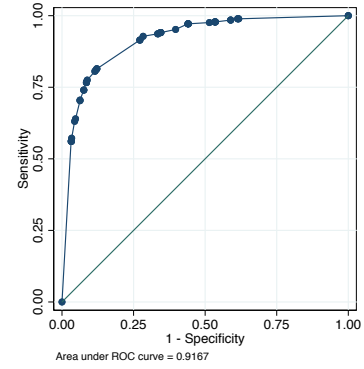
	Message Content	Promise?	Outcome
1	Considering some irrational guys, I will stop at C for get money not for angry	Yes	Share
2	If I were you I will stop at A. If you choose E, I would consider you as an irrational guys. You can be nice or not	No	Share
3	If you choose B, then I will choose C. Cross my heart and swear to die. :)	Yes	Reject
4	Have yourself a merry little Christmas	No	Accept
5	Hi, I hope you have a good rest of the day. Thanks for participating in this research.	No	Accept
6	hello love. Choose B and lets both make more money!!! i promise i'lll pick C	Yes	Reject
7	if you choose B, i'll choose C (you'll be making the same amount but also helping me—ut prosim amiright?)	Yes	Share
8	Hello, I am a poor, broke, college student, plz be reasonable and considerate and generous	No	Out
9	In my opinion, pineapple is a pretty good topping on pizza.	No	Accept
10	Thanks for choosing F	No	Out
11	Hokies play UVA in baseball today at 5:30 at home. Pick B and I'll choose C. Go Hokies!!	Yes	Share
12	Hello! I hope you're having a wonderful day :)	No	Reject
13	Don't YAWN if you yawn I am gonna pick C. If you do not YAWN I am probably gonna pick C... End of story, you may YAWN and I am gonna pick C	Yes	Share
14	If there was a meme or gif that could convey to you that I was picking C, I would send it to you 13 times not 17 because I want 13	Yes	Share
15	i know there is no way in hell we're making it to F so if you hook me up by picking B I will choose C	Yes	Share

C Beliefs, Plans, and Behavior

Supplementary Figure 1. P1's *Out* choices vs. P1's Plans

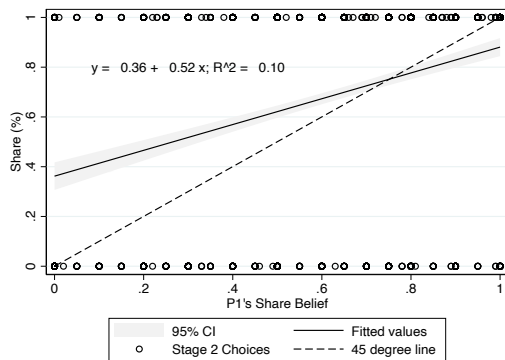


(a) Linear fit

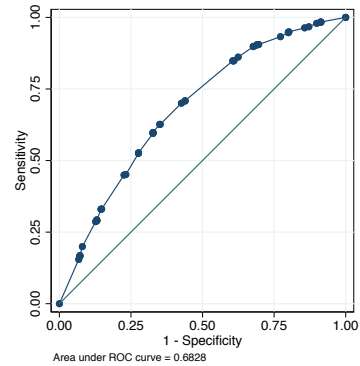


(b) ROC curve

Supplementary Figure 2. P2's *Share* choices vs. P1's Beliefs

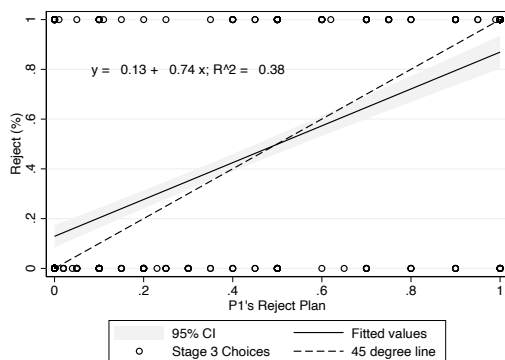


(a) Linear fit

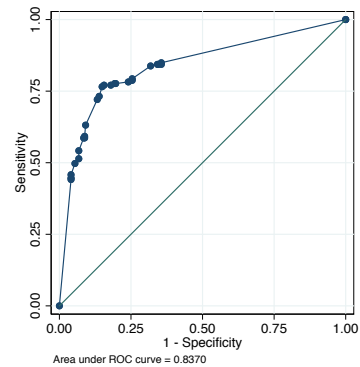


(b) ROC curve

Supplementary Figure 3. P1's *Reject* choices vs. P1's Plans

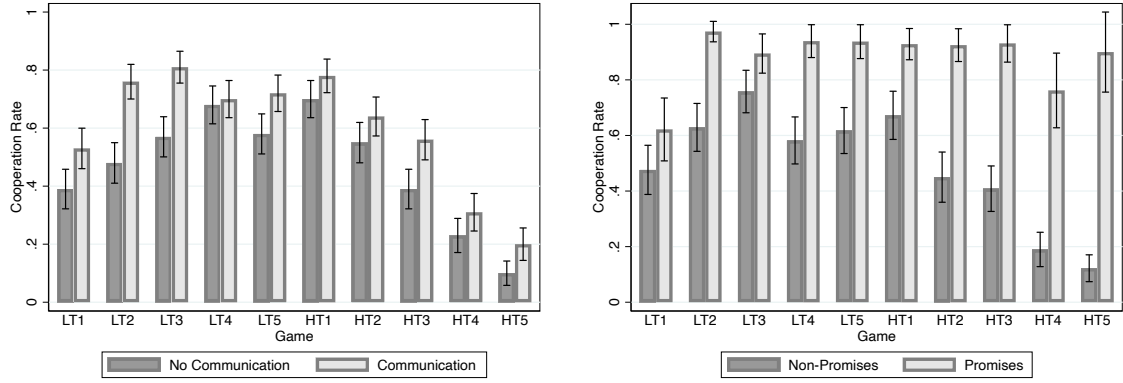


(a) Linear fit



(b) ROC curve

D Communication, Promises, and Cooperation



(a) High cooperation with communication.

(b) High cooperation with promises.

Supplementary Figure 4. Cooperation Rate by Communication and Promises.

References

- Aina, C., Battigalli, P., and Gamba, A. (2020). Frustration and anger in the ultimatum game: An experiment. *Games and Economic Behavior*, 122:150–167.
- Avoyan, A. and Ramos, J. (2020). A road to efficiency through communication and commitment. *SSRN 2777644*.
- Balliet, D. (2010). Communication and cooperation in social dilemmas: A meta-analytic review. *Journal of Conflict Resolution*, 54(1):39–57.
- Battigalli, P. and Dufwenberg, M. (2007). Guilt in games. *American Economic Review*, 97(2):170–176.
- Battigalli, P. and Dufwenberg, M. (2009). Dynamic psychological games. *Journal of Economic Theory*, 144(1):1–35.
- Battigalli, P. and Dufwenberg, M. (2022). Belief-dependent motivations and psychological game theory. *Journal of Economic Literature*, Forthcoming.
- Battigalli, P., Dufwenberg, M., and Smith, A. (2015). Frustration and anger in games. *IGIER 539*.
- Battigalli, P., Dufwenberg, M., and Smith, A. (2019). Frustration, aggression, and anger in leader-follower games. *Games and Economic Behavior*, 117:15–39.
- Berkowitz, L. (1989). Frustration-aggression hypothesis: Examination and reformulation. *Psychological Bulletin*, 106(1):59.
- Blanco, M., Engelmann, D., Koch, A. K., and Normann, H.-T. (2010). Belief elicitation in experiments: is there a hedging problem? *Experimental Economics*, 13(4):412–438.
- Blume, A. and Ortmann, A. (2007). The effects of costless pre-play communication: Experimental evidence from games with pareto-ranked equilibria. *Journal of Economic Theory*, 132(1):274–290.
- Bolton, G. E. and Ockenfels, A. (2000). ERC: A theory of equity, reciprocity, and competition. *American Economic Review*, pages 166–193.
- Bosman, R. and Van Winden, F. (2002). Emotional hazard in a power-to-take experiment. *The Economic Journal*, 112(476):147–169.

- Brandts, J. and Charness, G. (2011). The strategy versus the direct-response method: a first survey of experimental comparisons. *Experimental Economics*, 14(3):375–398.
- Cartwright, E. (2019). A survey of belief-based guilt aversion in trust and dictator games. *Journal of Economic Behavior & Organization*, 167:430–444.
- Charness, G. and Dufwenberg, M. (2006). Promises and partnership. *Econometrica*, 74(6):1579–1601.
- Che, Y.-K. and Sákovics, J. (2008). *Hold-Up Problem*. In: Palgrave Macmillan (eds) The New Palgrave Dictionary of Economics. Palgrave Macmillan, London.
- Crawford, V. P. (2016). New directions for modelling strategic behavior: Game-theoretic models of communication, coordination, and cooperation in economic relationships. *Journal of Economic Perspectives*, 30(4):131–50.
- Cuzick, J. (1985). A wilcoxon-type test for trend. *Statistics in medicine*, 4(1):87–90.
- Di Bartolomeo, G., Dufwenberg, M., and Papa, S. (2023). Promises and partner-switch. *Journal of the Economic Science Association*, 9(1):77–89.
- Dollard, J., Miller, N. E., Doob, L. W., Mowrer, O. H., and Sears, R. R. (1939). *Frustration and aggression*. Yale University Press.
- Dufwenberg, M. and Kirchsteiger, G. (2004). A theory of sequential reciprocity. *Games and Economic Behavior*, 47(2):268–298.
- Dufwenberg, M., Li, F., and Smith, A. (2025). Credible threats.
- Dufwenberg, M., Smith, A., and Van Essen, M. (2013). Hold-up: With a vengeance. *Economic Inquiry*, 51(1):896–908.
- Ellingsen, T. and Johannesson, M. (2004a). Is there a hold-up problem? *Scandinavian Journal of Economics*, 106(3):475–494.
- Ellingsen, T. and Johannesson, M. (2004b). Promises, threats and fairness. *The Economic Journal*, 114(495):397–420.
- Fehr, D. and Sutter, M. (2019). Gossip and the efficiency of interactions. *Games and Economic Behavior*, 113:448–460.
- Fehr, E. and Gächter, S. (2000). Fairness and retaliation: The economics of reciprocity. *Journal of Economic Perspectives*, 14(3):159–181.

- Fehr, E. and Schmidt, K. M. (1999). A theory of fairness, competition, and cooperation. *Quarterly Journal of Economics*, 114(3):817–868.
- Fischbacher, U. (2007). z-tree: Zurich toolbox for ready-made economic experiments. *Experimental Economics*, 10(2):171–178.
- Frank, R. H. (1988). *Passions within reason: the strategic role of the emotions*. WW Norton & Co.
- Geanakoplos, J., Pearce, D., and Stacchetti, E. (1989). Psychological games and sequential rationality. *Games and Economic Behavior*, 1(1):60–79.
- Grossman, S. J. and Hart, O. D. (1986). The costs and benefits of ownership: A theory of vertical and lateral integration. *Journal of Political Economy*, 94(4):691–719.
- Grout, P. A. (1984). Investment and wages in the absence of binding contracts: a Nash bargaining approach. *Econometrica*, pages 449–460.
- Hart, O. and Moore, J. (1990). Property rights and the nature of the firm. *Journal of Political Economy*, 98(6):1119–1158.
- Hart, O. and Moore, J. (2008). Contracts as reference points. *Quarterly Journal of Economics*, 123(1):1–48.
- Hirshleifer, J. (1987). On the emotions as guarantors of threats and promises. *The Dark Side of the Force*, pages 198–219.
- Klein, B., Crawford, R. G., and Alchian, A. A. (1978). Vertical integration, appropriable rents, and the competitive contracting process. *Journal of Law and Economics*, 21(2):297–326.
- Krupka, E. L., Leider, S., and Jiang, M. (2017). A meeting of the minds: Informal agreements and social norms. *Management Science*, 63(6):1708–1729.
- Murphy, K. P. (2012). *Machine learning: a probabilistic perspective*.
- North, D. C. and Weingast, B. R. (1989). Constitutions and commitment: the evolution of institutions governing public choice in seventeenth-century england. *Journal of Economic History*, 49(4):803–832.
- Persson, E. (2018). Testing the impact of frustration and anger when responsibility is low. *Journal of Economic Behavior & Organization*, 145:435–448.

- Rutström, E. E. and Wilcox, N. T. (2009). Stated beliefs versus inferred beliefs: A methodological inquiry and experimental test. *Games and Economic Behavior*, 67(2):616–632.
- Schotter, A. and Trevino, I. (2014). Belief elicitation in the laboratory. *Annual Review of Economics*, 6(1):103–128.
- Sell, A., Tooby, J., and Cosmides, L. (2009). Formidability and the logic of human anger. *Proceedings of the National Academy of Sciences*, 106(35):15073–15078.
- Selten, R. (1978). The chain store paradox. *Theory and Decision*, 9(2):127–159.
- Tirole, J. (1986). Procurement and renegotiation. *Journal of Political Economy*, 94(2):235–259.
- Toussaert, S. (2018). Eliciting temptation and self-control through menu choices: a lab experiment. *Econometrica*, 86(3):859–889.
- Trautmann, S. T. and Kuilen, G. (2015). Belief elicitation: A horse race among truth serums. *The Economic Journal*, 125(589):2116–2135.
- Vanberg, C. (2008). Why do people keep their promises? An experimental test of two explanations. *Econometrica*, 76(6):1467–1480.
- Williamson, O. E. (1971). The vertical integration of production: market failure considerations. *American Economic Review*, 61(2):112–123.
- Xiao, E. and Houser, D. (2011). Punish in public. *Journal of Public Economics*, 95(7-8):1006–1017.
- Yang, Y. (2021). A survey of the hold-up problem in the experimental economics literature. *Journal of Economic Surveys*, 35(1):227–249.