

Yale University

## EliScholar – A Digital Platform for Scholarly Publishing at Yale

---

Cowles Foundation Discussion Papers

Cowles Foundation

---

12-1-2013

### Promises and Expectations

Florian Ederer

Alexander Stremitzer

Follow this and additional works at: <https://elischolar.library.yale.edu/cowles-discussion-paper-series>



Part of the [Economics Commons](#)

---

#### Recommended Citation

Ederer, Florian and Stremitzer, Alexander, "Promises and Expectations" (2013). *Cowles Foundation Discussion Papers*. 2328.

<https://elischolar.library.yale.edu/cowles-discussion-paper-series/2328>

This Discussion Paper is brought to you for free and open access by the Cowles Foundation at EliScholar – A Digital Platform for Scholarly Publishing at Yale. It has been accepted for inclusion in Cowles Foundation Discussion Papers by an authorized administrator of EliScholar – A Digital Platform for Scholarly Publishing at Yale. For more information, please contact [elischolar@yale.edu](mailto:elischolar@yale.edu).

# PROMISES AND EXPECTATIONS

By

Florian Ederer and Alexander Stremitzer

December 2013  
Updated March 2016

COWLES FOUNDATION DISCUSSION PAPER NO. 1931



COWLES FOUNDATION FOR RESEARCH IN ECONOMICS  
YALE UNIVERSITY  
Box 208281  
New Haven, Connecticut 06520-8281

<http://cowles.yale.edu/>

# Promises and Expectations\*

Florian Ederer<sup>†</sup>  
Yale University

Alexander Stremitzer<sup>‡</sup>  
UCLA

March 14, 2016

## Abstract

We investigate why people keep their promises in the absence of external enforcement mechanisms and reputational effects. In a controlled laboratory experiment we show that exogenous variation of second-order expectations (promisors' expectations about promisees' expectations that the promise will be kept) leads to a significant change in promisor behavior. We provide clean evidence that a promisor's aversion to disappointing a promisee's expectation leads her to keep her promise. We propose a simple theory of lexicographic promise keeping that is supported by our results and nests the findings of previous contributions as special cases.

*Keywords:* promises, expectations, beliefs, contracts

*JEL Classification:* A13, C91, D03, C72, D64, K12.

---

\*We thank Mark Greenberg whose insightful comments influenced the design of our experiment. We are also grateful to Jason Abaluck, Jennifer Arlen, Pierpaolo Battigalli, Andreas Blume, Arthur Campbell, Gary Charness, Martin Dufwenberg, Christoph Engel, Florian Englmaier, Constança Esteves-Sorenson, Christine Exley, Robert Gibbons, Holger Herz, Lisa Kahn, Navin Kartik, Camelia Kuhnen, Rosario Macera Parra, Bentley MacLeod, Andreas Roider, Frédéric Schneider, Marta Serra Garcia, Seana Shiffrin, Joel Sobel, Rebecca Stone, Steven Tadelis, Noam Yuchtman, Kathryn Zeiler and seminar audiences at the Max Planck Institute for Research on Collective Goods, NYU, MIT, University of Regensburg, UCLA, Yale, ALEA, the NBER Summer Institute, the NBER Organizational Economics Meeting, the AEA 2015 meetings, SITE 2015, and the UCSD Deception Conference for helpful comments and suggestions as well as Estela Hopenhayn and the California Social Science Experimental Laboratory (CASSEL) at UCLA for helping us to conduct the experiments. We also thank Sean Maddocks, James Davis, and Jaimini Parekh for excellent research assistance and the UCLA Faculty Research Grant Program for financial support.

<sup>†</sup>Yale School of Management, 165 Whitney Avenue, New Haven, CT 06511, florian.ederer@yale.edu.

<sup>‡</sup>UCLA School of Law, 385 Charles E. Young Drive, 1242 Law Building, Los Angeles, CA 90095, stremitzer@law.ucla.edu.

# 1 Introduction

To facilitate production and exchange over time, parties often make promises in order to commit to a particular course of action. There are three main reasons why a party would honor such an obligation. The first is the existence of a third party enforcement mechanism.<sup>1</sup> A second reason stems from reputational concerns that arise when parties are concerned that renegeing on a promise might hurt their future payoffs.<sup>2</sup> A third reason, and the focus of the present paper, is the moral force of promise keeping.<sup>3</sup> A string of recent studies offers experimental evidence that promises, even if they come in the form of mere cheap talk, considerably enhance subsequent levels of cooperation in experimental hold-up, trust, and dictator games (Ellingsen and Johannesson 2004, Charness and Dufwenberg 2006, Vanberg 2008, Charness and Dufwenberg 2011).

While the practical relevance and effectiveness of the moral force of promise keeping is undisputed, there is a vigorous debate in economics, social psychology, philosophy, and law about why people keep (or should keep) their promises in the absence of explicit contractual and reputational concerns.<sup>4</sup> A clear understanding of what drives the desire to keep one's promise is essential to harnessing the beneficial effects of promises in institutional design, whether it be in the design of legal policy, regulatory regimes, contracts, or organizations.<sup>5</sup>

Two leading explanations for the moral force of promise keeping have been proposed. Proponents of the *expectation-based* theory argue that promisors (senders of promises) keep their word in order to avoid guilt when failing to meet the expectations the promise has created in promisees (receivers of promises). A promisor is therefore more likely to keep

---

<sup>1</sup>This is assumed in the economic literature on formal contracts beginning with Mirrlees (1976) and Holmström (1979).

<sup>2</sup>Such self-enforcing contracts have been studied extensively in the literature on relational contracting (Macaulay 1963, Klein and Leffler 1981, Bull 1987, Kreps 1990, MacLeod and Malcolmson 1989, Levin 2003).

<sup>3</sup>Note that in law every legally enforceable contract can at least partially rely on the moral force of promise keeping as a source of commitment. This is because the legal definition of a contract requires the existence of a promise as one of its elements (Restatement 2d of Contracts §1).

<sup>4</sup>Notable contributions to the broader literature on promise keeping in political sciences and social psychology include Ostrom, Walker, and Gardner (1992), Kerr and Kaufman-Gilliland (1994), Sally (1995), and Bicchieri and Lev-On (2007). In legal philosophy classic references include Fried (1981), Atiyah (1983), and Scanlon (1998). For a recent contribution containing a survey of the previous literature, see Shiffrin (2008).

<sup>5</sup>Finan and Schechter (2012) show that the moral force of promise keeping exists outside the lab and significantly influences political behavior. Even with secret balloting (and thus unobservable votes), promises to vote for a particular candidate are perceived to be sufficiently strong for vote-buying “candidates to target trustworthy voters who can be trusted to keep their promise” (p. 876). Eigen (2012) shows that reminding a contracting party that when entering into a contract they have made a promise is an effective way to assure compliance.

her promise if she believes that the promisee expects her to keep her promise.<sup>6</sup> In contrast, the *commitment-based* theory claims that promisors have a preference for keeping their word independent of the expectations of promisees. Promisors therefore suffer a cost from behaving in a way that is inconsistent with what they have promised.<sup>7</sup> The factors emphasized by these explanations are not mutually exclusive in explaining why people keep their promises. However, previous experimental research has either failed to disentangle these two motives or has only documented support for the commitment-based theory while suggesting that the expectation-based theory does not explain promise keeping. An implication of these results would be that making promisees' expectations salient in contractual and institutional design is not an effective way to induce more promise keeping or fulfillment of contractual obligations.

Using a novel design which exogenously varies expectations (while holding constant whether somebody makes a promise) this paper is the first to find clean experimental support for the expectation-based account. We propose a simple theory of *lexicographic* promise keeping in which a promisor is influenced by her promisee's expectation but only if these expectations are supported by the promise made by the promisor. This theory is consistent with our experimental results and nests the findings of previous contributions as special cases.

We use a trust game where a dictator (trustee) can make a free-form promise to a recipient (trustee) and the recipient can decide whether to trust the dictator and to remain in the game. Our main innovation is introduce a move of nature after this opt-in decision which determines the probability that it will be technically possible for the dictator to keep her promise. Both parties learn at this point whether the game is played with a "reliable random device" under which there is a high probability that the dictator will be able to keep her promise or whether the game is played with an "unreliable random device," under which there is a low probability that the dictator will be able to keep her promise. In the next step, another move of nature determines whether the dictator is able to perform or not. While both parties know with which random device the game is played, only the dictator

---

<sup>6</sup>Dufwenberg and Gneezy (2000) and Charness and Dufwenberg (2006) provide some suggestive experimental evidence in support of expectation-based theories.

<sup>7</sup>Experimental evidence for the commitment-based explanation for promise keeping can be found in the contributions of Braver (1995), Ostrom, Walker, and Gardner (1992), Ellingsen and Johannesson (2004), and Ismayilov and Potters (2012).

but not the recipient learns about whether she is able to perform or not. Therefore, a dictator who knows she is able to perform may face two kinds of recipients. Either the recipient has high expectations as he has learned that the game is played with the reliable random device, or, he has low expectations as he has learned that the game is played with the unreliable random device. This design allows us to compare promise keeping rates among dictators who are both able to keep their promises but hold different second-order beliefs (beliefs about how much the receiver expects to receive), depending on whether the history of the game leading up to the dictator’s decision reveals that it was likely (“reliable random device”) or unlikely (“unreliable random device”) that the dictator would be able to perform.

Using a within-subject design that allows us to observe dictators under both reliability settings, we show that the exogenous variation of the reliability of the random device with which the game is played directly affects the recipient’s first-order and the dictator’s second-order expectations, and that these significantly change the dictator’s decision to keep her promise. Our findings provide clean evidence for an expectation-based explanation of promise keeping: while the commitment created through promises between the two parties remains constant, second-order expectations increase due to the increase of the reliability of the random device which in turn induces an increase in promisors’ performance rates.

Finally, with a simple structural model we recover subject-specific susceptibilities to guilt aversion and characterize their distribution in the subject population. While slightly less than half of our subjects are unaffected by this behavioral trait, the remaining proportion exhibits some degree of guilt aversion and there is significant variation in how guilt-averse these subjects are.

**Related Literature** Our paper builds on the large literature on promises in economics, social psychology, philosophy, and law discussed in footnote 4. Methodologically, the paper is most closely related to the contributions of Charness and Dufwenberg (2006), Vanberg (2008), and Ellingsen et al (2010).

We show that promisees’ expectations play an important role when promisors decide whether to keep their promises in a trust game. However, this is only the case when there is a direct promissory link between the promisor and the promisee. In other words, dictator A’s performance decision is only influenced by recipient B’s expectations if A has made a promise to B. Where such a promissory link is missing, either because dictator A has made a

promise to recipient C, or dictator A has not made a promise at all, we find that recipients' expectations do not have a significant impact on dictators' performance choices. Yet, these latter cases without direct promissory links are exactly the scenarios considered in previous studies.

Using an ingenious experimental design, Vanberg (2008) argues that commitment *per se* is the only reason why people keep their promises. However, his suggestion that the expectation-based theory does not explain promise keeping is based on the observation that exogenous variation of second-order expectations does not lead to any statistically significant difference in promise keeping in a setting where there is no direct promissory link: When deciding whether to perform, a promisor is completely certain that she faces recipient B, and no longer recipient C to whom she initially made a promise.<sup>8</sup> Similarly, Ellingsen et al (2010) find no significant relationship between expectations and contributions in dictator/trust games in which they elicit expectations from recipients and communicate those beliefs to the dictator/trustee. Here too, a direct promissory link is missing because no promises are ever made. In contrast, in Charness and Dufwenberg (2006) there are direct promises and the authors provide suggestive evidence in support of guilt aversion, but their evidence, unlike ours, is based on correlations between second-order beliefs and actions that also admit alternative explanations.

The central contribution of our paper is to use exogenous variation in second-order beliefs and to show that the irrelevance of promisor expectations is a unique feature of the special cases in which expectations are not supported by direct promises. Instead, promise keeping seems to have a lexicographic structure: A promisee's expectations matter if and only if a direct promissory link exists between the promisor and the promisee.<sup>9</sup> Our theory and experiment can explain why guilt aversion plays a role in Charness and Dufwenberg (2006), but appears irrelevant in Vanberg (2008) and Ellingsen et al (2010). One reason for this

---

<sup>8</sup>Vanberg (2008) also studies situations in which the promissory link is not broken, but second-order expectations are not exogenously changed in that setting.

<sup>9</sup>Note that common law seems to track our lexicographic account of promise keeping. Under the doctrine of promissory estoppel, in the absence of an enforceable contract a disappointed party can recover against another party only if the other party made a promise that led to detrimental reliance (see Restatement (Second) of Contracts §90). More generally, across jurisdictions, law tends to protect reasonable expectations, but only if the injured party actually expected performance and only if the other side was in some way responsible for an act which created those expectations. By finding that the structure of promise keeping follows the same lexicographic structure, we find evidence consistent with a fundamental conjecture of legal theory: Law tracks and serves as a backstop to our moral intuitions.

lexicographic structure might be that a promise creates a sense of responsibility in the promisor, maybe because a promise establishes a personal connection which increases the salience of the promisee’s expectations or maybe because the promisor thinks that her act of promising *caused* the promisee’s expectations.

The remainder of the paper is organized as follows: Section 2 presents the design of the experiment, a simple model of promises and guilt aversion, and the experimental procedures. In Section 3 we report our results. Section 4 concludes. In Appendix A we report the exact instructions for the subjects participating in our experiment. Appendix B contains the formal proofs for our theoretical predictions.

## 2 Experimental Design and Procedure

We design an experiment to investigate the role of expectations in promise keeping. We hypothesize that a dictator is more likely to keep her promise if she believes that the recipient expects her to keep her promise. The underlying rationale for our hypothesis is that a dictator cares about the recipient’s expectations if and only if those expectations are supported by the dictator’s own promise.

Previous experiments were not intended to investigate this question, and hence they either confounded expectation- and commitment-based explanations or used the expectations created by other promisors as a means of varying the level of promisees’ expectations. Thus, our design is the first to shed light on what we refer to as the lexicographic structure of promise keeping. Instead of varying second-order beliefs through the promises given by a third party as in Vanberg (2008), the promissory link between the two parties is not broken in our experiment. Rather, the magnitude of the dictator’s expectations is exogenously varied by the type of random device that is selected.

### 2.1 Experimental Design

In our experiment, subjects are randomly matched in pairs in each period and play the experimental trust game depicted in Figure 1. In that game, after receiving a free-form message from the dictator, the recipient decides to opt in or out of the game and the dictator subsequently decides how much to contribute to the recipient.<sup>10</sup>

---

<sup>10</sup>This free-form message approach, which allows the dictator to send any message to the recipient (with the exception of identifying information such as name, age, race, gender) follows previous research by Charness



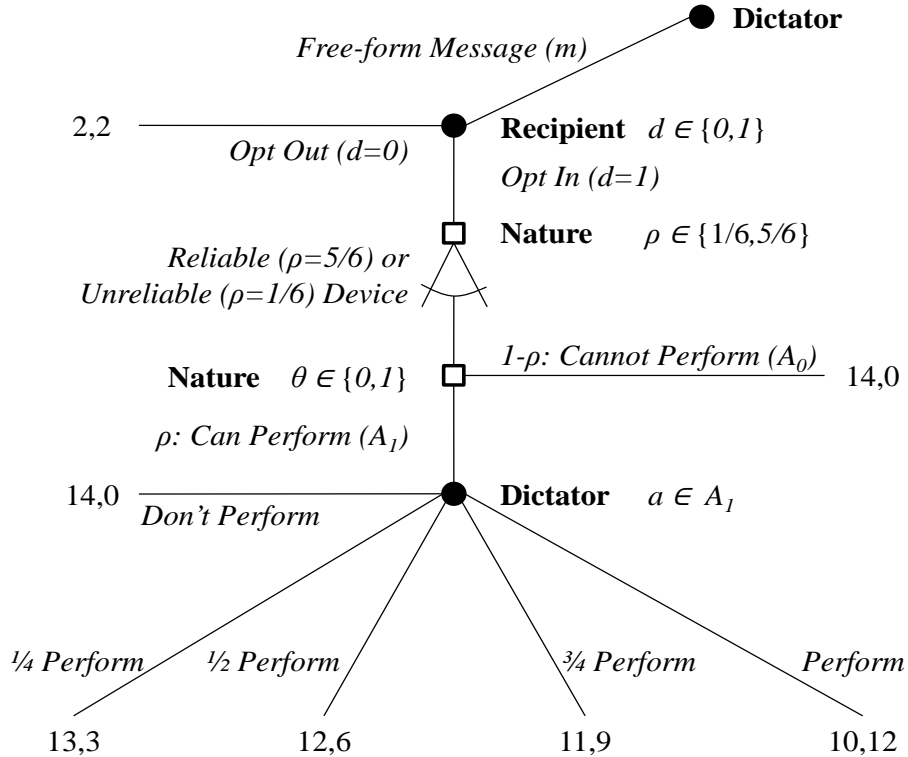


Figure 1: Dictator game with opt-out choice and reliable/unreliable device

The main feature of our design is that, following the opt-out decision of the recipient, nature selects whether the subjects play the game with a *reliable* or an *unreliable* device. This device determines how likely it is that the dictator will be able to choose some positive level of performance (i.e., any action other than *Don't Perform*). If the random device is reliable, the dictator can choose an action that delivers a positive payoff to the recipient (*Perform*, *3/4 Perform*, *1/2 Perform*, *1/4 Perform*) with probability  $5/6$ . If the random device is unreliable, the dictator can only deliver on a promise to perform with probability  $1/6$ . For example, if performance is impossible, the dictator receives \$14 and the recipient receives \$0. If performance is possible and the dictator chooses *Perform*, she receives \$10 and the recipient receives \$12. Figure 1 depicts the remaining payoffs for the two players. If the dictator chooses *Don't Perform* the parties receive the same payoffs (14, 0) as if performance had not been possible.

Formally speaking, the random device determines how likely it is that the dictator finds and Dufwenberg (2006) and Vanberg (2008). In contrast, Charness and Dufwenberg (2010) use pre-coded messages and find only small effects of such ‘bare’ promises.

herself in one of two states of the world,  $\theta \in \{0, 1\}$ , with associated action space  $A_\theta$ . This action space depends on whether the dictator is able to perform,  $\theta = 1$ , or not  $\theta = 0$ :

$$a \in A_\theta = \begin{cases} A_0 = \{0\} & \text{if } \theta = 0 \\ A_1 = \{0, .25, .5, .75, 1\} & \text{if } \theta = 1 \end{cases} .$$

Figure 1 illustrates that the dictator's monetary payoffs can be written as a function of her decision  $a$  by  $\pi_D(a) = 14 - 4a$ . Similarly the recipient's monetary payoffs are represented by  $\pi_R(a) = 12a$ .

The timing of the game is as follows. At the beginning of each period ( $t = 0$ ), subjects are randomly paired and nature randomly determines the identity of the second mover (dictator, promisor) and the first mover (recipient, promisee) in each pair. At  $t = 1$  the dictator can send a free-form message  $m$  that, following the literature on promises, is coded as no message, empty talk, or promise,  $\mu(m) \in \{\emptyset, 0, 1\}$ . Although we distinguish between the different messages in our experimental data, for the purpose of our model, we will treat no message and empty talk as the same category such that  $\mu(m) \in \{0, 1\}$ . At  $t = 2$ , the recipient decides to opt out or to stay in,  $d \in \{0, 1\}$ . If  $d = 0$ , the game ends and payoffs (2, 2) for the dictator and the recipient are realized. If  $d = 1$ , the game continues to  $t = 3$ . At  $t = 3$ , nature determines the type of the random device  $\rho$  with which the game is played. The random device can be either *reliable* ( $\rho = 5/6$ ) or *unreliable* ( $\rho = 1/6$ ) where  $\rho$  denotes the probability that the dictator will be able to choose from action space  $A_1$ . Both parties learn the type of the random device. At  $t = 4$ , the dictator but not the recipient learns the state of the world  $\theta$  and she makes the decision  $a \in A_\theta$ . Thus the dictator knows that she can perform but, when making her choice, she faces recipients who plausibly expect a higher payoff under the reliable than under the unreliable scenario. At  $t = 5$ , both players learn their payoffs, and the recipient learns the state  $\theta$ .<sup>11</sup>

In addition to recording the dictator's choice  $a$  at  $t = 4$  in our experiment, we elicit the recipient's and the dictator's beliefs at  $t = 3$ . The recipient is asked which action  $a \in A_1$  the dictator is going to choose if she will be able to perform at  $t = 4$ . Because the recipient knows

---

<sup>11</sup>This design feature means that if the dictator chose not to perform although she was able to, the recipient will know that the dictator broke her promise. The dictator cannot hide behind the circumstances. Eliminating this deniability is important as a dictator would benefit from a higher level of deniability in the unreliable ( $\rho = 1/6$ ) as opposed to the reliable ( $\rho = 5/6$ ). As we explain later, our main design idea is to use the two random devices to exogenously vary second-order expectations and it is important that we do not introduce a confound by having the identity of the random device also influence the level of deniability. Tadelis (2011) theoretically and experimentally shows that the degree of exposure of the dictator's actions to the recipient's scrutiny also influences the dictator's actions.

with which random device  $\rho$  the game is played, these beliefs might depend on the history of the game (realizations of  $\mu$  and  $\rho$ ). We therefore denote the first-order belief of the recipient by  $\tau_R(\mu, \rho) \in [0, 1]$ . In turn, the dictator has a belief (probability measure) regarding  $\tau_R$ . Let  $\tau_D(\mu, \rho) \in [0, 1]$  denote the mean second-order belief the dictator has about the recipient’s belief  $\tau_R(\mu, \rho)$ , which we also elicit at  $t = 3$ .<sup>12</sup> As we mention in our detailed description of the experimental procedure, this elicitation of beliefs was incentivized.<sup>13</sup> It is important to note at  $t = 4$  (i.e., at the time when the dictator makes her decision) the dictator—but not the recipient—knows the state of  $\theta$  (i.e., whether the dictator can perform).

Except for the slightly richer action set for the dictator and a within-subject design, this game is to a large extent identical to the trust/dictator game in Charness and Dufwenberg (2006) and Vanberg (2008).<sup>14</sup> We chose a richer action space to allow for more variation in the contribution rates of dictators as well as a within-subject design that asks dictators to choose actions for both the reliable and unreliable device to increase statistical power.<sup>15</sup> However, what really distinguishes our design from the previous two papers is the existence of the random device, which determines the probability  $\rho$  with which the dictator will be able to choose some positive level of performance. The main purpose of this design innovation is to exogenously vary the dictator’s and the recipient’s expectations without breaking any promissory link that exists between a dictator and a receiver.<sup>16</sup>

If the random device is reliable, then there is a probability of  $\rho = 5/6$  that the dictator will be able to choose *Perform*. If, on the other hand, the random device is unreliable, there is only a probability of  $\rho = 1/6$  that the dictator will be able to choose *Perform*. Thus, recipients

---

<sup>12</sup>Similar to the design in Vanberg (2008) the elicitation of the beliefs in our experiment restricts the set of first- and second-order beliefs to a set of five elements that mirrors  $A_1$ .

<sup>13</sup>The main benefit of incentivizing beliefs is that it increases the accuracy of beliefs. As a downside, there is evidence that eliciting beliefs and, in particular, incentivizing belief elicitation affects contribution levels in repeated public goods games, although the direction of the effect seems very sensitive to design specifics (e.g., Croson 2000, Gächter and Renner 2010). This does not pose a problem for our design, given that we are only interested in the difference between contribution levels between the reliable and unreliable scenario which should cancel out any level effects. We should be more concerned about an increased steepness between beliefs and contributions as our design is based on manipulating second-order beliefs. However, there is only weak evidence of a very small effect of belief elicitation on the relationship between beliefs and contribution rates (Gächter 2010).

<sup>14</sup>Vanberg (2008) also uses euros rather than dollars, thus generating slightly higher payoffs for the participants.

<sup>15</sup>Brandts and Charness (2000) have shown that decisions elicited using the strategy method do not differ significantly from those elicited using the “hot” method.

<sup>16</sup>As discussed before, in Vanberg (2008), recipient expectations are exogenously varied by using third-party promises, but this variation comes at the expense of a broken promissory link between the two parties. In Ellingsen et al (2010) no promises are ever made to begin with.

who are playing the game with an unreliable random device can plausibly expect lower monetary payoffs from the game. Because dictators are aware of this change in expectations (due to independent variation in the experiment), their second-order expectations are also exogenously changed. It is important to note that our manipulations cannot affect the commitment *per se* because at the time the promise is made, the dictator only knows that the game is potentially played with the reliable or the unreliable random device, but not which of the two scenarios will subsequently be realized. Similarly, at the time the dictator makes her decision at  $t = 4$ , she—but not the recipient—knows whether she is able to perform independent of which random device the game was played with. At this point, only the history of the game differs. If higher second-order expectations lead to higher contribution rates by the dictators who promised to perform, this would constitute evidence for the expectation-based explanation of promise keeping.<sup>17</sup>

To state our ideas more formally, we now present a simple model that builds on psychological game theory (Geanakoplos, Pearce and Stacchetti 1989, Battigalli and Dufwenberg 2007, 2009) and captures the effect of guilt aversion on promise keeping.

## 2.2 A Simple Model of Promises and Guilt Aversion

Battigalli and Dufwenberg (2007) propose two general theories of guilt aversion based on simple guilt and guilt from blame, respectively. Our model uses the former concept of simple guilt in which the dictator cares about the extent to which she lets the recipient down.<sup>18</sup> We define  $\gamma_D \geq 0$  as a constant measuring the dictator’s sensitivity to guilt from disappointing the recipient’s expectations, which the dictator expects to be equal to  $E[\pi_R | \tau_D, \rho] = \rho 12 \tau_D(\rho)$ . We conjecture that promise keeping has a lexicographic structure. Recipients’ expectations only play a role if the dictator has made a promise ( $\mu = 1$ ), but not if she has not made a promise ( $\mu = 0$ ) to the recipient.

The dictator’s utility  $U_D$  when she chooses  $a$  at  $t = 4$  can now be written in the following

---

<sup>17</sup>It is important to remember that in  $t = 5$  the recipient learns  $\theta$ , that is, whether the dictator was able to perform or not. See footnote 11.

<sup>18</sup>A philosopher might object to this terminology as it comes close to depicting guilt aversion as a primary moral motivation. It can only be a secondary one, activated by a person’s belief that some other fact itself provides direct motivation to act. In this case, the prospect of disappointing expectations would be wrong and therefore one should not do it.

way:

$$\begin{aligned}
U_D(a) &= \pi_D(a) - \frac{\mu\gamma_D}{k} (\max\{E[\pi_R|\tau_D(\mu, \rho), \rho] - \pi_R(a), 0\})^k \\
&= 14 - 4a - 12^k \frac{\mu\gamma_D}{k} (\max\{\rho\tau_D(\mu, \rho) - a, 0\})^k.
\end{aligned} \tag{1}$$

The last term of the dictator’s utility function captures the impact of guilt. This term only plays a role if the dictator sent a promise ( $\mu = 1$ ) and if the dictator is susceptible to guilt aversion ( $\gamma_D > 0$ ). Guilt from disappointing the recipient’s perceived expectations  $E[\pi_R|\tau_D, \rho]$  by choosing a low payoff  $\pi_R(a)$  for the recipient has a negative effect on utility, but there is no gain from exceeding the recipient’s expectations. The dictator can reduce the negative utility from guilt by increasing her action  $a$  up to the point where it matches the dictator’s beliefs about the recipient’s expectations. In contrast to Charness and Dufwenberg (2006) and Battigalli and Dufwenberg (2007, 2009), we allow guilt to be linear ( $k = 1$ ) or convex ( $k > 1$ ) in the difference between the dictator’s expectations,  $E[\pi_R|\tau_D, \rho]$ , and the realized payoffs for the recipient,  $\pi_R(a)$ . For  $k = 1$ , our model nests the model of Charness and Dufwenberg (2006) as a special case that only admits corner solutions of  $a$ . For  $k > 1$ , interior solutions of  $a$  (i.e.,  $a \neq \{0, 1\}$ ) are also possible.<sup>19</sup> Finally, and most importantly, note the role of the random device. When the random device is reliable ( $\rho = 5/6$ ), the impact of guilt—by virtue of the larger expected payoff for the recipient—is much larger than when the device is unreliable ( $\rho = 1/6$ ).<sup>20</sup>

### 2.2.1 Predictions

It is instructive to establish two benchmark cases in which expectations do not affect actions. First, a dictator who is motivated solely by her own monetary payoff and is not sensitive to guilt at all,  $\gamma_D = 0$ , would have a utility of  $14 - 4a$  and would therefore maximize her payoff by choosing  $a = 0$ .<sup>21</sup> Second, in the settings considered by Vanberg (2008) and Ellingsen et al (2010), in which no direct promissory link between the dictator and the recipient exists and

---

<sup>19</sup>For ease of presentation we focus on the latter case in the main text. In the appendix, we show that our results also hold for  $k = 1$ .

<sup>20</sup>Our model draws on the theory of simple guilt aversion developed by Battigalli and Dufwenberg (2009) who adapt the sequential equilibrium concept of Kreps and Wilson (1982) to psychological games. However, for the most important part of our theoretical analysis (**TP1** below), like Charness and Dufwenberg (2006), we do not invoke any equilibrium supposition.

<sup>21</sup>Of course, there are many reasons other than guilt aversion such as social preferences or norms (Rabin 1993, Fehr and Schmidt 1999, Bolton and Ockenfels 2000, Andreoni and Bernheim 2009) that would predict an equilibrium action  $a$  other than 0.

thus  $\mu = 0$ , beliefs about expectations  $\tau_D(\rho)$  also do not matter and hence the predictions are the same as in our first benchmark.

Thus, our model requires two assumptions for second-order beliefs to play a role in promise keeping. There must exist a promise between the two parties,  $\mu = 1$ , and the dictator must experience some guilt aversion,  $\gamma_D > 0$ . Second-order expectations will then generate different predictions about the contribution choice  $a$  for the reliable ( $\rho = 5/6$ ) and the unreliable ( $\rho = 1/6$ ) device. The dictator's utility is given by

$$U_D = 14 - 4a - 12^k \frac{\gamma_D}{k} (\max\{\rho\tau_D(\rho) - a, 0\})^k \quad (2)$$

which yields different levels of guilt for the two different devices and where we write  $\tau_D(\mu = 1, \rho) = \tau_D(\rho)$  for simplicity. It is straightforward to see that the impact of guilt is larger for  $\rho = 5/6$  than for  $\rho = 1/6$  thus leading to a higher equilibrium action  $a$  for two reasons. First, there is a difference in actions resulting purely from the exogenous variation in the reliability  $\rho$  of the device. Second, there is an additional (second-order) effect resulting from the impact of this exogenous variation on equilibrium beliefs  $\tau_D$ . The first-order condition with respect to  $a$  for the dictator yields the following interior solution:

$$a^* = \rho\tau_D(\rho) - \left(\frac{4}{\gamma_D 12^k}\right)^{\frac{1}{k-1}} \quad (3)$$

The dictator's action  $a^*$  is increasing in the reliability of the random device,  $\rho$ , the dictator's second-order belief,  $\tau_D$ , and her susceptibility to guilt aversion,  $\gamma_D$ .

To see the first effect of  $\rho$  on  $a$ , assume that second-order beliefs about actions are the same in both settings,  $\tau_D(5/6) = \tau_D(1/6)$ , and that, just for ease of exposition,  $k = 1$ . As can be seen from equation (2), the guilt experienced by the dictator when choosing  $a = 0$  is  $2\gamma_D\tau_D$  for  $\rho = 1/6$ , which is much smaller in magnitude relative to the guilt experienced for  $\rho = 5/6$  where it is  $10\gamma_D\tau_D$ . This argument holds a fortiori for  $\tau_D(5/6) > \tau_D(1/6)$  as is evident from the first-order condition for interior solutions of  $a$  from equation (3). Thus, in equilibrium, we expect the dictator to choose higher levels of  $a$  for  $\rho = 5/6$  than for  $\rho = 1/6$ .<sup>22</sup> This leads us to the main hypothesis generated by our model.

**Theoretical Prediction 1** *If there is a promise ( $\mu = 1$ ), the dictator's contribution action  $a$  is higher for the reliable than for the unreliable device. If there is no promise ( $\mu = 0$ ), there is no difference in the dictator's contribution action. (TP1)*

---

<sup>22</sup>See Appendix B for a more rigorous proof of the theoretical predictions taking the possibility of corner solutions into account.

As a result of these different action choices, we also expect first- and second-order beliefs to differ in the two settings. In particular, because predicted actions are higher for  $\rho = 5/6$ , equilibrium first-order beliefs of recipients must adjust to the different actions that the dictator chooses in the two settings. Hence, first-order beliefs are higher,  $\tau_R(5/6) \geq \tau_R(1/6)$  and, as a result, equilibrium second-order beliefs must be higher too,  $\tau_D(5/6) \geq \tau_D(1/6)$  if  $\gamma \geq 0$ . This leads us to our second prediction.

**Theoretical Prediction 2** *If there is a promise ( $\mu = 1$ ), first-order and second-order beliefs are higher for the reliable device than for the unreliable device. If there is no promise ( $\mu = 0$ ), there is no difference in first-order and second-order beliefs. (TP2)*

At this point it is crucial to remind the reader that in contrast to the central prediction about contribution actions (TP1) this prediction regarding beliefs (TP2) is a little more subtle because it requires beliefs to adjust to the changes in actions postulated in TP1. That is to say, in our model, any differences in beliefs are caused by differences in contribution actions between the two random devices and message categories. If, however, in terms of contribution actions dictators do not respond differently to the random devices and messages then beliefs will also be unaffected. Note further that this difference in second-order beliefs provides a second reason for why actions are higher under  $\rho = 5/6$  and reinforces the first effect on the dictator's action discussed in TP1. Because the dictator's action is higher in equilibrium if there is a promise ( $\mu = 1$ ) compared to when there is none ( $\mu = 0$ ), the recipient's expected payoff from not opting out is higher if  $\mu = 1$ . Therefore, the recipient's opt-out decision and his beliefs should be responsive to the message of the dictator. This leads us to our third hypothesis.

**Theoretical Prediction 3** *The recipient's opt-in rate and first-order beliefs are higher if there is a promise ( $\mu = 1$ ) than if there is no promise ( $\mu = 0$ ). (TP3)*

Consistent with the results obtained by Charness and Dufwenberg (2006) and Vanberg (2008), we find that promises made in the communication phase are correlated with a significant increase in both the rate at which dictators choose an action other than *Don't Perform* and in the second-order beliefs concerning the probability that they will do so. However, the primary goal of the present paper is to investigate whether independent variation in the

level of second-order beliefs, achieved through the use of a random device, leads to different performance levels by dictators who previously made promises.<sup>23</sup>

## 2.3 Experimental Procedure

We conducted 20 experimental sessions with a total of 280 student subjects at the California Social Science Experimental Laboratory (CASSEL). The CASSEL subject pool consists of undergraduate students from UCLA. Subjects were assigned to visually isolated computer terminals. Beside each terminal they found paper instructions, which are reproduced in Appendix A. Questions were answered individually at the subjects' seats.

Each session consisted of 2 unpaid practice rounds followed by 8 paying rounds. In each round, subjects interacted with another randomly chosen participant. Under no circumstances did any participant interact with the same participant twice in the paying rounds. We achieved this by creating matching groups of exactly 10 subjects. At the end of the experiment, one of the 8 paying rounds was randomly chosen for payment. Each round was equally likely to be selected. The amount paid out at the end of the experiment depended on the decisions made in that round. In each period we also elicited first- and second-order beliefs of subjects about the behavior of other subjects. This elicitation of beliefs was incentivized and subjects were paid for all rounds except the one chosen for payment of the decision to prevent hedging. The subjects received a fixed fee of \$10 for arriving on time. The experiment was programmed and conducted with the software *z-Tree* (Fischbacher 2007).

First, each subject was randomly matched with an interaction partner, and one participant in each pair was randomly assigned to the role of a dictator or the role of a recipient.<sup>24</sup> The pair matches and the roles of the players were randomly assigned anew in each round. It was always equally likely to be assigned to either role, regardless of the previous messages or actions in the game.

Second, each dictator could send free-form messages to the recipient. The dictator could

---

<sup>23</sup>We do not (intend to) show that the fact of having made a promise *causes* performance rates to increase. It might be that dictators who make a promise are just more likely to perform, independent of whether they have made a promise or not. To rule out such a selection effect, we would have to run a treatment in which it is not possible to make a promise and compare the performance rates with those chosen in a treatment where making a promise is possible. This is exactly the main manipulation in Charness and Dufwenberg (2006) who find evidence supporting a causal inference.

<sup>24</sup>In the instructions, we neutrally refer to the role of the dictator and the role of the recipient as "Role A" and "Role B," respectively.



send any number of (unidirectional) messages of a length of 256 characters each within a time frame of 90 seconds.<sup>25</sup> Subjects were not allowed to reveal their identity to the other participant. That is, they were not allowed to reveal their name or any other identifying feature such as race, gender, hair color, or seat number. In every other respect, subjects were free to send any message they liked.

Third, after receiving the message of the dictators, the recipients could decide whether to opt out. If a recipient chose to opt out, each player received \$2. If a recipient chose not to opt out, the game continued. At this point, neither player knew whether the random device determining whether the dictator would be able to perform was *Random Device 5/6* (“reliable random device”, probability of 5/6 that dictator would be able to choose something other than *Don’t Perform*) or *Random Device 1/6* (“unreliable random device”, probability of 1/6 that dictator would be able to choose something other than *Don’t Perform*). However, both parties knew that each scenario could occur with an equal probability of 50%.

Fourth, nature privately determined whether the players would play the game with the reliable or the unreliable random device, which determined the probability with which the dictator in the pair was able to choose a positive level of performance. At this point neither player learned with which type of random device the game was played. The players only learned at the end of each round which random device was chosen.

Fifth, the recipient guessed which choice the dictator would likely make if she could choose to perform and the dictator guessed which payoff the recipient expected to earn. Specifically, recipients and dictators were asked to choose from a five point scale. While the recipient’s guessing payoff depended on the contribution decision of the dictator, the dictator’s guessing payoff depended on the belief chosen by the recipient. Both payoffs were higher the closer they were to the actual contribution and belief, respectively.<sup>26</sup> As both parties did not know which random device had been chosen, we asked the players to make their guesses for both reliability scenarios. Note that if a recipient thought that the dictator intended to choose *Perform* (allocating \$12) a recipient’s expected payoff depended on the reliability scenario. The expected payoff was \$2 if the game was played with Random Device 1/6 ( $12 \times 1/6 = 2$ ) and \$10 if the game was played with Random Device 5/6 ( $12 \times 5/6 = 10$ ). Asking the

---

<sup>25</sup>Note that the 90-second time frame was not enforced, as it just served as an informal pacemaker.

<sup>26</sup>With the exception of five (rather than two) potential choices for the contribution decision our belief elicitation method is identical to that used in Vanberg (2008).

dictator to make her guesses in terms of the recipient’s expected payoffs allowed us to make those expectations particularly salient. These procedures yielded five point scales between 0 (performance very unlikely) and 1 (performance very likely) for first- and second-order beliefs.

Sixth, the dictator was asked to assume that she were able to perform and to make their contribution decisions as if the game leading up to that point had been played with Random Device 5/6 or Random Device 1/6. Figure 1 depicts the players’ payoffs under the different possible contribution decisions.

Seventh, the computer randomly drew an equally likely integer between 1 and 6 for each pair using z-Tree’s random number generator. If the random device was reliable it was possible for the dictator to perform for the numbers 2, 3, 4, 5, and 6. If the random device was unreliable, the dictator was able to perform for the number 1.

Finally, at the end of each round, both dictators and recipients learned with which random device the game had been played, whether it had been possible for the dictators to perform, and the payoffs both participants had earned.

### 3 Results

The data comprise 20 experimental sessions involving a total of 280 subjects with a total of 28 matching groups of 10 subjects. Each session lasted for 8 rounds. The average number of dictator decisions made by each subject is 4. As we used the strategy method to elicit first- and second-order beliefs and contribution choices for both types of the random device, this within-subject design gives us a total of 1,120 decisions made under each reliability scenario. However, each matching group constitutes only one independent observation. Non-parametric tests are therefore based on matching group averages of the relevant variables. For comparisons between the random devices, we have matched observations that allow us to use two-tailed Wilcoxon signed-rank tests while for unmatched comparisons between message categories we use two-tailed Mann-Whitney ranksum (MW) tests.<sup>27</sup>

---

<sup>27</sup>For all the instances where the MW ranksum test is used, we also used the Fligner-Policello (FP) test. This test is a robust rank test for unmatched data which does not require that the two populations that are to be compared have equal variances (see Section 4.4 of Hollander et al (2014) for a complete description of the FP test). In our analysis the p-values of each MW ranksum test and its corresponding FP test are almost identical and we therefore omit reporting the results for the latter test.

### 3.1 Performance Rates

To investigate the role of promises, we asked a student assistant to code messages according to whether they contained a promise stating that the subject would choose any action other than *Don't Perform*. Following Charness and Dufwenberg (2006), this classification yielded three categories: “no message” ( $\mu = \emptyset$ ) containing no text at all; “empty talk” ( $\mu = 0$ ) messages (e.g., “Hey I just met you/and this is crazy/but here’s my message/so money maybe?”); and “promise” ( $\mu = 1$ ) messages (e.g., “im going to choose 3/4th perform so please dont opt out”). After accounting for all opt-out decisions (see Section 3.4), there remain 383, 300, and 268 individual observations, and 28, 27, and 28 matching group observations in the three message categories (promise, empty talk, no message) for both the reliable and the unreliable device. We use these for our remaining analysis.

When the dictator made a promise ( $\mu = 1$ ), the average contribution (performance rate) she gave to the recipient was \$7.08 (0.59) for the reliable ( $\rho = 5/6$ ) and \$6.48 (0.54) for the unreliable device ( $\rho = 1/6$ ), conditional on performance being feasible. While this difference is small in magnitude it is highly statistically significant (Wilcoxon signed-rank,  $p$ -value  $< 0.01$ ) and consistent with our central hypothesis **TP1**. When the dictator made a promise, she actually chose higher performance rates when it subsequently turned out that it was likely that she would be able to perform (reliable device), as opposed to when the possibility of performance was unlikely (unreliable device). Furthermore, as we show below, the small magnitude of the contribution differentials is in line with the small, but significant differentials in second-order beliefs between the two random device settings, lending further empirical support to the expectations-based explanation for promise keeping. Dictators contributed significantly more under the reliable device scenario because they (correctly) held exogenously higher second-order beliefs about the recipients’ expectations in that scenario.

As we are employing the strategy method (and thus a within-subject design), it is particularly instructive to examine the behavior of those dictators who made different contribution decisions in the unreliable and the reliable device setting. Figure 2 shows the contribution decisions of dictators who promised to contribute and who chose to alter their contribution decision depending on whether the device was reliable or unreliable.<sup>28</sup> A much lower proportion of dictators chose *Don't Perform* for the reliable than for the unreliable device, and

---

<sup>28</sup>Note that all of our statistical tests are based on the full sample of dictators and not just those dictators who changed their decision between the two random devices.

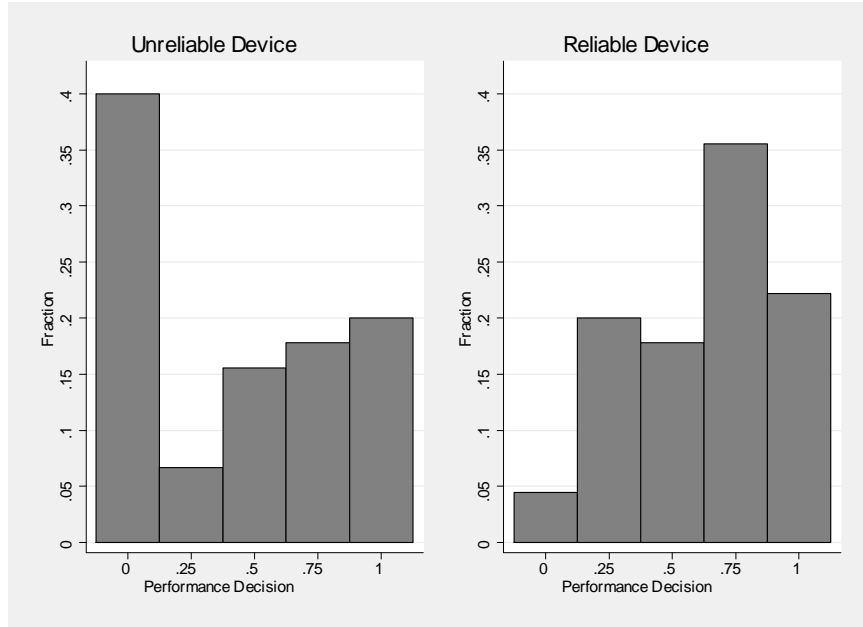


Figure 2: Fraction of performance decisions of dictators (*Don't Perform*,  $1/4$  Perform,  $1/2$  Perform,  $3/4$  Perform, Perform) who sent a promise and changed their decision between the unreliable and the reliable device

hence more of them ended up choosing higher performance rates. While 40% of dictators chose *Don't Perform* under the unreliable device, only less than 5% chose the same action under the reliable device. As a result, a much larger proportion of dictators chose to contribute a positive amount under the reliable device, with the  $3/4$  Perform action experiencing the largest increase.<sup>29</sup>

At this point, it is crucial to note that the significant difference in dictator contributions between the two settings disappeared when the dictator engaged in empty talk (Wilcoxon signed-rank,  $p$ -value 0.27) or sent no message at all ( $p$ -value 0.14). Thus, in the absence of an explicit promise, higher reliability did not lead to higher performance. This finding is the empirical equivalent of the second part of **TP1**, which shows that contributions should not differ between the two reliability settings if the dictator did not promise. This evidence is also consistent with the conjectured lexicographic structure of promise keeping. The receiver's expectations influence the dictator's contribution decisions only when the dictator committed herself to a promise; they play little or no role when the dictator made no promise.

<sup>29</sup>Figure 2 also shows that many subjects make interior choices of  $a$ , suggesting that guilt aversion is not linear ( $k = 1$ ), but that it is instead convex ( $k > 1$ ) in the difference between beliefs and actions. We explore the magnitude of guilt aversion in greater detail in Section 3.5.

We want to stress that while our design provides crisp evidence for the central thesis of our paper that recipients' expectations matter for dictators' performance decisions, our results are only consistent with the conjectured lexicographic structure of promise keeping. Our design is not ideally suited for directly testing this lexicographic structure. The fact that expectations do not matter for those dictators who did not make a promise could be due to a selection effect: dictators who are more likely to promise are also more likely to be affected by differences in recipients' expectations. However, some support for our conjecture could be derived from the fact that Vanberg (2008) found no effect of recipients' expectations on the dictators' decision to perform when there was no promissory link, while we clearly identify a positive effect on performance when such a direct promissory link exists.

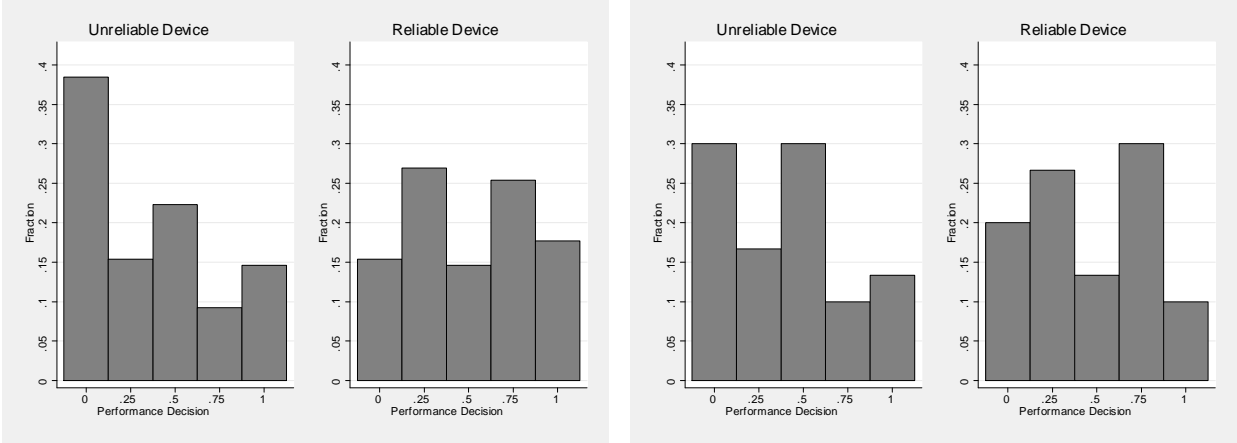


Figure 3a: Same as Figure 2 for empty talk    Figure 3b: Same as Figure 2 for no message

The two panels of Figure 3 depict the same data as Figure 2 for those dictators who changed their contribution decision between the two settings, but instead focuses on the message categories of empty talk and no message. While there is a small change in behavior towards more generous contribution rates from the unreliable to the reliable device for dictators who sent an empty talk message (Figure 3a), there is practically no change for dictators who sent no message (Figure 3b). The slightly positive (but statistically insignificant) shift in dictator contribution rates is larger for empty talk messages than for no messages. In a related experiment that investigates commitment-based explanations of promise keeping Ismayilov and Potters (2012) find that both trustees who make a promise *and* those who do not are more likely to be trustworthy if their message is delivered to the trustor. Their

findings as well as ours suggest that any form of communication increases trustworthiness irrespective of the content of messages. This tentatively suggests that the guilt aversion effect is larger when there is some communication rather than none at all. However, as pointed out above, in both cases, and in contrast to when the dictator made a promise, there is no significant difference in the mean of the contribution decisions.

Finally, note that the highly statistical difference in performance rates between the reliable and unreliable device for promises (Wilcoxon signed-rank,  $p$ -value  $< 0.01$ ) and the lack of statistical significance for empty talk ( $p$ -value 0.27) and no message ( $p$ -value 0.14) are not caused by large differences in the number of individual observations (and hence statistical power) across the three categories. As mentioned above, there are 383 (promise), 300 (empty talk), and 268 (no message) individual observations, and 28, 27, and 28 matching group observations in the three message categories for both the reliable and the unreliable device. In accordance with other previous contributions, we also find that average performance rates are higher if the dictator made a promise (0.56) than if there was just empty talk (0.38) or no message (0.38) (MW rank-sum,  $p$ -values  $< 0.01$ ). This translates into a contribution difference of \$2.16 which is about 3 times as large as the \$0.60 difference in contributions between the two random devices though that difference is purely due to differences to second-order expectations. In contrast, the \$2.16 difference is a combination of several effects. It contains a selection effect (subjects who promise are different from those who do not promise), the commitment effect (subjects feel compelled to contribute just because of the promise *per se*) and the expectation effect. Furthermore, in contrast to face-to-face interactions, between people who know each other well and may care strongly about not disappointing the other party's expectations, we would not expect the expectation effect to be very large in magnitude in our anonymous, low interaction experimental laboratory setting. However, we show that even in this setting which is very similar to consumer-to-consumer eCommerce transactions, expectations matter.

### 3.2 First-Order Beliefs and Expectations

Having documented evidence for the expectation-based explanation by analyzing the differences in contribution rates across treatments, we now investigate whether the secondary predictions of our model regarding beliefs are also borne out in our data. Recipients were asked to guess the dictator's decision on a five-point scale between 0 and 1. When the

random device was reliable, the recipients had mean first-order beliefs  $\tau_R(\mu, 5/6)$  of 0.23, 0.26, and 0.34 for no message, empty talk and promise, respectively. These beliefs were slightly lower (0.22, 0.21, and 0.29) for the same messages in the case of an unreliable device,  $\tau_R(\mu, 1/6)$ . For both outcome realizations of the random device, receiving a promise (relative to receiving no message or empty talk) significantly raised the recipients' expectations about how much they would receive from the dictator, moving first-order expectations from the lower values of 0.23 and 0.26 to 0.34 and from 0.22 and 0.21 to 0.29, respectively (MW rank-sum,  $p$ -values 0.007, 0.005, 0.022, 0.006). This pattern mirrors the results of Charness and Dufwenberg (2006) and Vanberg (2008) and is consistent with **TP3** of our model.

However, although the first-order beliefs are lower across the board for the unreliable device, the differences in first-order beliefs between the two device settings are not statistically significant (Wilcoxon signed-rank,  $p$ -values 0.17, 0.95, 0.45). Even when focusing on the case in which the dictator sent a promise, the difference in the recipient's first-order beliefs is not statistically significant between a reliable and an unreliable device. The lack of a statistically significant difference in first-order beliefs between the two reliability settings is consistent with **TP2** if the dictator did not promise (no message and empty talk). However, it is inconsistent with the same theoretical prediction for the case where the dictator made a promise, because the model predicts a difference in first-order beliefs in that case.

It is important to remember that we elicited conditional beliefs that directly match  $\tau_R(\mu, \rho)$  in our model and allow for easy comparability across the two reliability devices. Recipients were asked how much (on a five-point scale) they thought the dictator would contribute if performance was feasible, as the recipients did not actually learn whether performance was feasible for the dictator until after the end of each round. Therefore, the relevant first-order *expectations* at the time the dictator forms her second-order expectations and makes her performance decision are given by the unconditional first-order expectations, which are  $12\tau_R(\rho)\rho$ . Thus, in order to obtain the first-order expectations in terms of expected payoffs, the elicited conditional first-order beliefs have to be multiplied by  $5/6 \times 12 = 10$  for the reliable and by  $1/6 \times 12 = 2$  for the unreliable device. We find that these unconditional first-order beliefs are substantially higher in the reliable (\$2.60, \$2.30, and \$3.40) than in the unreliable scenario (\$0.42, \$0.44, and \$0.58), and that these differences are statistically significant (Wilcoxon signed-rank,  $p$ -values  $< 0.01$ ).<sup>30</sup>

---

<sup>30</sup>This stark difference is a (mechanical) feature of our experimental design. In contrast to previous

### 3.3 Second-Order Beliefs and Expectations

We next investigate how second-order beliefs  $\tau_D(\mu, \rho)$  (i.e., a dictator’s belief  $\tau_D(\mu, \rho)$  about the belief  $\tau_R(\mu, \rho)$  that the recipient has about the dictator’s performance decision) vary with the message sent by the dictator and the reliability of the random device. When the random device was reliable, second-order beliefs were 0.59 when the dictator sent no message, 0.59 for an empty talk message, and 0.72 when a promise was given. In contrast, when the random device was unreliable, the same second-order beliefs fell to 0.48, 0.43, and 0.64, respectively. The second-order beliefs in both reliability settings are significantly higher for promises than for empty talk and no messages (MW rank-sum,  $p$ -values  $< 0.01$ ).

More importantly, though, given the focus of this paper, if the dictator made a promise the difference in second-order beliefs between the reliable and the unreliable random device is statistically significant (Wilcoxon signed-rank,  $p$ -value  $< 0.01$ ) which is consistent with **TP2**. That is to say, when a dictator made a promise to share money with the recipient, her belief about the amount the recipient expected she would contribute was significantly higher when the random device was reliable (0.72) than when it was unreliable (0.64). The likely reason for this finding is the dictator’s realization that the recipient expected a higher level of performance when the random device was reliable than when it was unreliable which led the dictator to adapt her second-order beliefs accordingly.

When the dictator did not send any message or sent an empty talk message, second-order beliefs were also significantly higher in the 5/6 random device case than in the 1/6 random device condition (Wilcoxon signed-rank,  $p$ -values  $< 0.01$ ). This contradicts **TP2** which predicts that second-order beliefs should not differ between the two reliability settings if the dictator did not make a promise. However, as shown in Section 3.1 and as predicted by our main hypothesis **TP1** contribution rates were not significantly higher when expectations were not supported by a promise.<sup>31</sup>

---

contributions, we do not vary second-order expectations through the endogenous variation of first-order beliefs. Instead, we directly and exogenously change second-order expectations through the different random device scenarios and the timing of when the dictator and the recipient learn about which random device was chosen.

<sup>31</sup>Dictators may have thought that even if they gave no promise recipients may hold higher (out-of-equilibrium) first-order expectations if the random device is reliable than if it is unreliable. Although this contradicts **TP2**, this out-of-equilibrium phenomenon may lend additional support to our theory of lexicographic promise keeping. Although the second-order beliefs are higher they do not translate into higher contribution rates if there is no promise. In other words, differences in second-order beliefs only matter if they are justified by a promise.



The second-order expectations given by  $E[\pi_R|\tau_D, \rho] = 12\tau_D(\rho)\rho$  that correspond to these second-order beliefs influence the level of guilt experienced by the dictator subjects as can be seen from the utility function in equation (1). These second-order expectations are equal to \$5.88, \$5.91, and \$7.18 for the reliable device and significantly lower (Wilcoxon signed-rank,  $p$ -values  $< 0.01$ ) for the unreliable device where they are equal to \$0.96, \$0.87, and \$1.28. This large difference in second-order expectations is, of course, mainly exogenously (and mechanically) created by the use of the random device and this large variation serves us in our goal to separately identify the effect of guilt aversion on promise keeping.

### 3.4 Promises and Opt-Out Decisions

Recipients chose to opt out at different rates depending on which message they received from the dictator with whom they were paired. While only 7.3% of recipients chose to opt out after receiving a message classified as a promise, 21.6% opted out if they received no message at all, and 12.8% opted out if they received an empty talk message. These differential opt-out rates are consistent with **TP3**, which predicts that the recipient’s (expected) payoff from staying in the game is higher if the dictator sent a promise than if she did not.

The differences in opt-out rates between recipients who received a promise and those who received an empty talk message, and between empty talk messages and no messages are only significant at the 10% level (MW rank-sum,  $p$ -values 0.07, 0.07). In contrast, the difference in opt-out rates between participants who received a promise and those who did not receive any message is significant (MW rank-sum,  $p$ -value  $< 0.01$ ). Similarly, the difference between participants who received a promise and (pooled) participants who did not receive a promise, either because they received an empty talk message or no message at all, is significant at the 1% level (MW rank-sum,  $p$ -value  $< 0.01$ ). The low opt-out rate of recipients who received a promise from their partnered dictator indicates that the recipients expected higher relative payoffs from staying in the game than from opting out compared to recipients who received no message at all or just an empty talk message. Furthermore, the higher opt-out rate for recipients who received no message relative to recipients of an empty talk message suggests that some form of verbal engagement is better than none at all when it comes to inducing recipients to stay in the game.<sup>32</sup>

---

<sup>32</sup>This is in line with the aforementioned findings of Potters and Ismayilov (2012) who also find that even some form of communication like empty talk increases trustworthiness.

### 3.5 Distribution of $\gamma_D$

Using our simple model of guilt aversion as well the experimental data obtained on performance actions  $a$ , beliefs  $\tau_R$  and  $\tau_D$ , and message categories  $\mu$ , we can directly recover each dictator's susceptibility to guilt aversion,  $\gamma_D$ . However, because many dictators choose strictly positive levels of performance  $a$  even when there is no promise ( $\mu = 0$ ), we augment our previous model by an additional term that captures altruism. The dictator's utility function is then given by

$$U_D = \pi_D(a) - \mu \frac{\gamma_D}{k} (\max\{E[\pi_R|\tau_D, \rho] - \pi_R(a), 0\})^k - \frac{\delta_D}{r} (12 - \pi_R(a))^r$$

where the dictator suffers a convex disutility if the receiver's payoff  $\pi_R(a)$  falls short of his maximum possible payoff of 12.<sup>33</sup> If  $\delta_D = 0$ , the dictator is not driven by altruism, but as  $\delta_D$  increases she cares more about the payoff obtained by the receiver. The values of  $k$  and  $r$  influence the convexity of the guilt aversion and the altruism terms, and are assumed to be known by the dictator and the recipient. Using identifying variation for subjects who are observed in the data choosing  $a$  under both a promise ( $\mu = 1$ ) and no promise ( $\mu = 0$ ), the unknown altruism and guilt aversion parameters,  $\delta_D$  and  $\gamma_D$ , are exactly identified from the first-order conditions with respect to  $a$ . To see this ignore, just for expositional purposes, corner solutions and note that for  $\mu = 0$  we have the following first-order condition for the dictator's choice of  $a$ :

$$\frac{\partial U_D}{\partial a} = -4 + 12\delta_D (12 - 12a)^{r-1} = 0.$$

This allows us to obtain  $\delta_D$ . Similarly, for  $\mu = 1$  we have

$$\frac{\partial U_D}{\partial a} = -4 + 12\gamma_D (\max\{12\rho\tau_D - 12a, 0\})^{k-1} + 12\delta_D (12 - 12a)^{r-1} = 0$$

which using  $\delta_D$  yields  $\gamma_D$ .<sup>34</sup>

The two panels of Figure 4 show the distribution of the altruism and guilt aversion parameters in the dictator population for quadratic altruism,  $r = 2$ , and quadratic guilt

<sup>33</sup>By using altruism we chose a very simple form of other-regarding preferences. This is because more general specifications could potentially lead to higher-order beliefs also playing a role. For example, dictators could be making contributions according to what is expected of them independent of their initial statement.

<sup>34</sup>In our experimental setting dictators are constrained to choose  $a \in \{0, .25, .5, .75, 1\}$ . However, since we also have several observations of  $a$  for each dictator we use dictator-specific averages of  $a$  for given  $\mu$  and  $\rho$  when solving the first-order conditions for  $\delta_D$  and  $\gamma_D$ . We take corner solutions at  $a = \{0, 1\}$  into account.

aversion,  $k = 2$ .<sup>35</sup> Given our assumptions, we are only able to identify the distribution of  $\gamma_D$  for 104 dictators who are observed under both  $\mu = \{0, \emptyset\}$  and  $\mu = 1$ . As suggested by our reduced form analysis that documents a significant positive shift in performance from unreliable to reliable device when  $\mu = 1$ , the distribution of  $\gamma_D$  (Figure 4b) shows that more than half of the dictators exhibit guilt aversion,  $\gamma_D > 0$ , while the remaining, slightly smaller proportion of just under 50% is unaffected by this behavioral trait,  $\gamma_D = 0$ .

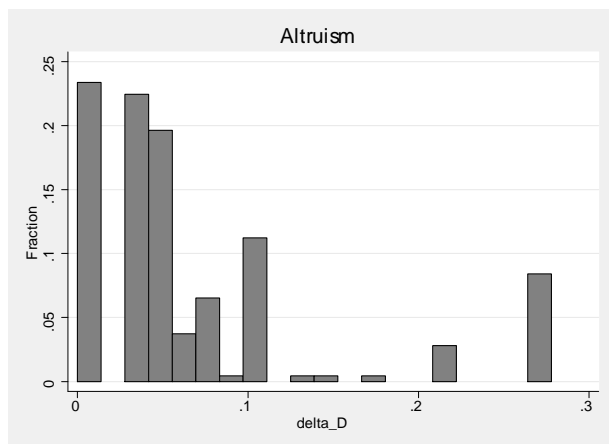


Figure 4a: Altruism parameter  $\delta_D$  for  $r = k = 2$

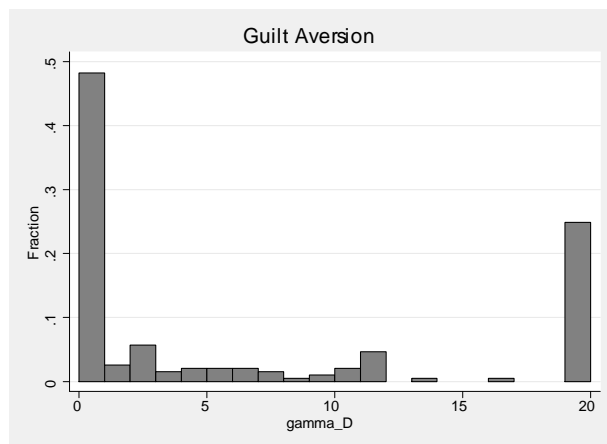


Figure 4b: Guilt aversion parameter  $\gamma_D$  for  $r = k = 2$

The dictators with a positive  $\gamma_D$  fall into two broad categories as can be seen from Figure 4b. First, there is a mass of about 25% of all dictators (at  $\gamma_D \approx 20$ ) where  $\gamma_D$  is so large that in equilibrium the dictators raise  $a$  sufficiently high such that  $a \approx \rho\tau_D$ . In this way, they reduce their own monetary payoff in order to completely avoid any loss from guilt aversion which they would otherwise suffer if they chose a lower performance  $a$ . Of course, the true  $\gamma_D$  of these dictators might be even higher than 20, but these subjects are already at a corner solution in our data. Second, for roughly 25% of dictators,  $\gamma_D$  lies between the two mass points of 0 and 20, and so these dictators trade off some monetary gains against losses from guilt aversion. They do not, however, raise  $a$  high enough to completely eliminate guilt in equilibrium.

<sup>35</sup>Note that while convexity in the guilt aversion term (i.e.,  $k > 1$ ) is required to explain any subject choices that are not corner solutions, we chose these particular parameters for their simplicity and their fit with our data. Our results are qualitatively unchanged when we used different parameters.

## 4 Conclusions

Many psychological and economic experiments have shown that promises greatly enhance cooperative behavior in experimental games, but evidence on the driving forces behind why people keep their promises remains scant. In this paper we provided the first clean evidence for the expectation-based explanation of promise keeping.

Previous experiments either could not distinguish between commitment-based and expectation-based explanations, because treatment-induced changes in the alternative causal factors (promises and second-order beliefs) had occurred simultaneously, or focus on settings in which there was an insufficient level of commitment between the dictator and the recipient for the changes in the levels of (first- and second-order) expectations to affect behavior.

In contrast, we designed our experiment to achieve independent variation in second-order expectations in an environment where these were supported by a direct promissory link between the dictator and the recipient, and thus by the existence of a sufficiently high level of commitment. Changes in the probability with which a dictator would be able to contribute directly impacted recipients' first-order and dictators' second-order expectations, which in turn significantly changed behavior. In light of our own as well as previous findings, we propose a lexicographic theory of promise keeping. Guilt and expectations matter if and only if they are supported by the promise between the two acting parties.

Strictly speaking, Vanberg (2008) shows two things: First, he demonstrates that expectations *per se*, that is, expectations unsupported by a direct promise, do not explain promise keeping. In other words, he shows that without a direct promise expectations are insufficient to induce cooperation. Second, he shows that promises have a positive effect on cooperation when there are positive expectations. The present paper shows that expectations have a positive effect on cooperation if there is a promise. However, evidence that promising *per se* is sufficient to induce cooperation is so far missing. To see this, imagine that the only motivation for keeping a promise is that the promisor does not want to disappoint the promisee's expectations. Further assume that, consistent with our lexicographic theory, the sensitivity to the recipient's expectations is only switched on by a promise. Then, under the expectation-based theory, we would predict that a dictator who has made a promise is more likely to perform than a dictator who has not made a promise simply because her sensitivity to expectations will only be switched on if she made a promise. Note that for this to hap-

pen we need *not* assume any independent preference for promise keeping. The *whole* effect could work exclusively through the desire not to disappoint expectations so long as they are supported by a promise. In order to show that there is an independent preference for promise keeping we would have to design an experiment that varies the dictator's promissory commitment while keeping the recipient's expectations at zero. This would be an interesting avenue for future research.

## References

- ANDREONI, J., AND B. D. BERNHEIM (2009): "Social Image and the 50-50 Norm: A Theoretical and Experimental Analysis of Audience Effects," *Econometrica*, 77(5), 1607–1636.
- ATIYAH, P. (1983): *Promises, Morals, and Law*. Clarendon Press.
- BATTIGALLI, P., AND M. DUFWENBERG (2007): "Guilt in Games," *American Economic Review, Papers and Proceedings*, 97(2), 170–176.
- (2009): "Dynamic Psychological Games," *Journal of Economic Theory*, 144(1), 1–35.
- BICCHIERI, C., AND A. LEV-ON (2007): "Computer-Mediated Communication and Cooperation in Social Dilemmas: An Experimental Analysis," *Politics, Philosophy, and Economics*, 6(2), 139–168.
- BOLTON, G. E., AND A. OCKENFELS (2000): "ERC: A Theory of Equity, Reciprocity, and Competition," *American Economic Review*, 90(1), 166–193.
- BRANDTS, J., AND G. CHARNES (2000): "Hot vs. Cold: Sequential Responses and Preference Stability in Experimental Games," *Experimental Economics*, 2(3), 227–238.
- BRAVER, S. L. (1995): "Social Contracts and the Provision of Public Goods," in *Social Dilemmas: Perspectives on Individuals and Groups*, ed. by D. Schroeder. Praeger, New York.
- BULL, C. (1987): "The Existence of Self-Enforcing Implicit Contracts," *Quarterly Journal of Economics*, 102(1), 147–159.
- CHARNESS, G., AND M. DUFWENBERG (2006): "Promises and Partnership," *Econometrica*, 74, 1579–1601.
- (2010): "Bare promises: An experiment," *Economics Letters*, 107(2), 281–283.
- (2011): "Participation," *American Economic Review*, 101(4), 1211–1237.

- CROSON, R. T. A. (2000): “Thinking Like a Game Theorist: Factors Affecting the Frequency of Equilibrium Play,” *Journal of Economic Behavior & Organization*, 41(3), 299–314.
- DUFWENBERG, M., AND U. GNEEZY (2000): “Measuring Beliefs in an Experimental Lost Wallet Game,” *Games and Economic Behavior*, 30(2), 163–182.
- EIGEN, Z. (2012): “When and Why Individuals Obey Contracts: Experimental Evidence of Consent, Compliance, Promise, and Performance,” *Journal of Legal Studies*, 41, 67–93.
- ELLINGSEN, T., AND M. JOHANNESSON (2004): “Promises, Threats and Fairness,” *Economic Journal*, 114(495), 397–420.
- ELLINGSEN, T., M. JOHANNESSON, S. TJØTTA, AND G. TORSVIK (2010): “Testing Guilt Aversion,” *Games and Economic Behavior*, 68(1), 95–107.
- FEHR, E., AND K. M. SCHMIDT (1999): “A theory of fairness, competition, and cooperation,” *Quarterly Journal of Economics*, 114(3), 817–868.
- FINAN, F., AND L. SCHECHTER (2012): “Vote-Buying and Reciprocity,” *Econometrica*, 80(2), 863–881.
- FISCHBACHER, U. (2007): “z-Tree: Zurich Toolbox for Ready-Made Economic Experiments,” *Experimental Economics*, 10, 171–178.
- FRIED, C. (1981): *Contract as Promise*. Cambridge MA, Harvard University Press.
- GÄCHTER, S., AND E. RENNER (2010): “The Effects of (Incentivized) Belief Elicitation in Public Goods Experiments,” *Experimental Economics*, 13(3), 364–377.
- GEANAKOPOLOS, J., D. PEARCE, AND E. STACCHETTI (1989): “Psychological Games and Sequential Rationality,” *Games and Economic Behavior*, 1(1), 60–79.
- HOLLANDER, M., D. WOLFE, AND E. CHICKEN (2013): *Nonparametric Statistical Methods*, Wiley Series in Probability and Statistics. Wiley.
- HOLMSTRÖM, B. (1979): “Moral Hazard and Observability,” *Bell Journal of Economics*, 10, 74–91.
- ISMAYILOV, H., AND J. POTTERS (2012): “Promises as Commitments,” *Tilburg University, Center for Economic Research, Discussion Paper*, 2012-064.
- KERR, N. L., AND C. M. KAUFMAN-GILLILAND (1994): “Communication, Commitment and Cooperation in Social Dilemma,” *Journal of Personality and Social Psychology*, 66(3), 513–529.
- KLEIN, B., AND K. B. LEFFLER (1981): “The Role of Market Forces in Assuring Contractual Performance,” *Journal of Political Economy*, 89(4), 615–641.
- KREPS, D. M. (1990): “Corporate Culture and Economic Theory,” in *Perspectives on Positive Political Economy*, ed. by J. E. Alt, and K. A. Shepsle, pp. 90–143. Cambridge University Press.

- KREPS, D. M., AND R. WILSON (1982): “Sequential Equilibria,” *Econometrica*, 50(4), 863–894.
- LEVIN, J. (2003): “Relational Incentive Contracts,” *American Economic Review*, 93(3), 835–857.
- MACAULAY, S. (1963): “Non-Contractual Relations in Business: A Preliminary Study,” *American Sociological Review*, 28, 55–69.
- MACLEOD, B., AND J. MALCOMSON (1989): “Implicit Contracts, Incentive Compatibility, and Involuntary Unemployment,” *Econometrica*, 57, 447–480.
- MIRRELES, J. (1976): “The Optimal Structure of Incentives and Authority Within an Organization,” *Bell Journal of Economics*, 7, 105–131.
- OSTROM, E., J. WALKER, AND R. GARDNER (1992): “Covenants With and Without a Sword: Self-Governance Is Possible,” *American Political Science Review*, 86(2), 404–417.
- RABIN, M. (1993): “Incorporating Fairness into Game Theory and Economics,” *American Economic Review*, 83(5), 1281–1302.
- SALLY, D. (1995): “Conversation and Cooperation in Social Dilemmas: A Meta-Analysis of Experiments from 1958 to 1992,” *Rationality and Society*, 7(1), 58–92.
- SCANLON, T. (1998): *What We Owe to Each Other*. Cambridge MA, Harvard University Press.
- SHIFFRIN, S. V. (2008): “Promising, Intimate Relationships, and Conventionalism,” *Philosophical Review*, 117, 481–524.
- TADELIS, S. (2011): “The Power of Shame and the Rationality of Trust,” *UC Berkeley Haas Working Paper*.
- VANBERG, C. (2008): “Why Do People Keep Their Promises? An Experimental Test of Two Explanations,” *Econometrica*, 76, 467–1480.

## A Instructions

Thank you for participating in this experiment. The purpose of this experiment is to study how people make decisions in a particular situation. In case you should have questions at any time, please raise your hand. Please do not speak to other participants during the experiment. You will receive \$10 for arriving on time. Depending on the decisions made and the decisions of other participants, you may receive an additional amount (as described below). At the end of the experiment, the entire amount will be paid to you individually and privately in cash.

This session consists of 2 practice rounds and 8 paying rounds with money prizes. In each round, you will interact with another randomly chosen participant. Under no circumstances will you interact with the same participant twice. No participant will learn the identity of the persons with whom he or she has interacted during the experiment.

At the end of the experiment, one of the 8 paying rounds will be randomly chosen for payment (every round is equally likely). The amount that you will receive at the end of the experiment will depend on the decisions made in that round.

Each round consists of 7 steps, which are described below.

### Overview

There are two players; Player A and Player B. Initially, A can send a chat message to B over the computer, and B can decide whether he wants to opt out of the game, leading to payoffs of \$2 for each player. If B does not opt out, a random device will determine whether it will be possible for A to perform, that is, allocate money to B. If it is impossible to perform, Player A gets a payoff of \$14 and B gets a payoff of 0. If it is possible for A to perform, he can make one of 5 choices:

- *Don't Perform*: A keeps \$14 for himself and allocates \$0 to B.
- $1/4$  *Perform*: A keeps \$13 for himself and allocates \$3 to B.
- $1/2$  *Perform*: A keeps \$12 for himself and allocates \$6 to B.
- $3/4$  *Perform*: A keeps \$11 for himself and allocates \$9 to B.
- *Perform*: A keeps \$10 for himself and allocates \$12 to B.

There are two types of Random Device

- Random Device  $5/6$ : A is able to choose something other than *Don't Perform* with probability  $5/6$ .
- Random Device  $1/6$ : A is able to choose something other than *Don't Perform* with probability  $1/6$ .

The players learn about the type of the random device after B has made his opt-out decision.

**Step 1: Role assignment.** At the beginning of each round, you will be anonymously and randomly matched with another participant. Each member of the pair will then be randomly assigned Role A or Role B with equal probability (50%).

**Step 2: Communication.** During the communication phase, Player A can send a chat message to Player B. Important: You are not allowed to reveal your identity to the other participant. (That is, you may not reveal your name or any other identifying feature such as race, gender, hair color, or seat number.) In every other respect, you are free to send any message you like. Please continue to remain quiet while communicating with the other participant. Participants who violate these rules (experimenter discretion) will be excluded from the experiment and all payments.

**Step 3: Opt-out decision.** Player B can decide whether to opt out. If B chooses to opt out, each player receives \$2. If B chooses not to opt out, the game continues. Information: Neither player knows, whether the Random Device determining if A will be able to choose



*Perform* is Random Device 5/6 (probability that A can choose something other than *Don't Perform* is 5/6) or Random Device 1/6 (probability that A can choose something other than *Don't Perform* is 1/6). However, both parties know that each scenario occurs with equal probability (50%).

**Step 4: Nature of the Random Device revealed.** The players learn whether they play with Random Device 5/6 or Random Device 1/6.

**Step 5: Guessing.** Player B guesses which choice Player A is likely to make in Step 7. A guesses which payoff B expects to gain. Note that if B thinks that A intends to choose *Perform*—allocating \$12—B's expected payoff depends on what B has learned about the Random Device: The expected payoff is \$2 if the game is played with Random Device 1/6 ( $12 \times 1/6 = 2$ ) and \$10 if the game is played with Random Device 5/6 ( $12 \times 5/6 = 10$ ).

**Step 6:** Player A learns whether he will be able to perform. If only *Don't Perform* is possible, the game ends. If A is able to perform, the game continues to Step 7.

**Step 7: Decision phase.** A decides whether to choose *Don't Perform* (keep \$14 and send \$0 to B), or whether to choose *Perform* (keep \$10 and send \$12 to B) or any of the options in between. The payoffs are

	A	B
A chooses <i>Don't Perform</i>	\$14	\$0
A chooses $1/4$ <i>Perform</i>	\$13	\$3
A chooses $1/2$ <i>Perform</i>	\$12	\$6
A chooses $3/4$ <i>Perform</i>	\$11	\$9
A chooses <i>Perform</i>	\$10	\$12
B chooses "Opt Out"	\$2	\$2
Performance not possible	\$14	0

**Information at the end of a round.** Players learn their own payoff, which random device was chosen, and the players learn whether player A was able to perform.

**Conditional Choice.** You will be asked to make the guess in Step 5 and the decision in Step 7 before Step 4 has actually been played. In other words, you will be asked to assume that A will be able to perform in Step 7, and then make the guess in Step 5 and the decision in Step 7 for two scenarios:

1. Random device 1/6 was chosen.
2. Random device 5/6 was chosen.

Subsequently, Steps 4 and 6 are played and A's recorded choice will be entered as A's decision in Step 7 (provided the game reaches this step). A's decision will influence payoffs as if A took the same decision in Step 7.

**Bonus: Guessing.** At certain points, you will have the additional possibility to earn a small amount by guessing the decisions of the other participant. Guessing will be paid in every round that is not chosen for payment of the decision. You will learn more about this during the experiment.

**Do you have any questions?**

## B Proofs

We use Battigalli and Dufwenberg's (2007) general model of simple guilt to capture guilt aversion in our model. Applying their formulation to our game and notation yields the following utility function for the dictator

$$U_D = \pi_D(a) - \gamma_D \max \{E[\pi_R|\tau_D, \rho] - \pi_R(a), 0\}.$$

To this formulation we add the lexicographic structure of promise keeping which is governed by  $\mu$  in our model, and we allow for convex guilt to obtain

$$U_D = \pi_D(a) - \frac{\mu\gamma_D}{k} (\max \{E[\pi_R|\tau_D, \rho] - \pi_R(a), 0\})^k$$

We will first prove **TP1** by showing that the dictator's equilibrium choice  $a^*$  after having given a promise is strictly increasing in  $\rho$  for sufficiently high guilt aversion and 0 otherwise. The prediction for  $\mu = 0$  follows trivially from the dictator's utility function in expression (2). We distinguish two cases,  $k = 1$  and  $k > 1$ .

**Case  $k = 1$ .** If  $k = 1$ ,  $U(a)$  is a linear function in  $a$  and  $U'(a) = -4 + 12\gamma_D$ . It follows that the equilibrium action  $a^*$  maximizing the dictator's utility is given by the following corner solutions:

$$a^* = \begin{cases} 0 & \text{if } \gamma_D \leq \frac{1}{3} \\ \rho\tau_D & \text{if } \gamma_D > \frac{1}{3}. \end{cases} \quad (4)$$

**TP1** for  $k = 1$  follows directly from (4).

**Case  $k > 1$ .** If  $k > 1$ , note that  $U'(a) = -4 < 0$  on the interval  $[\rho\tau_D, \infty)$ . Therefore, the only candidate  $\hat{a}_1$  for a maximizer of the dictator's utility function on the interval  $[\rho\tau_D, \infty)$  is the corner solution  $\rho\tau_D$ , which is increasing in  $\rho$  and  $\tau_D$ .

On the interval  $[0, \rho\tau_D)$ , note that the dictator's utility function is strictly concave as  $U''(a) = -12^k(k-1)\gamma_D(\rho\tau_D - a)^{k-2} < 0$  for all  $a \in [0, \rho\tau_D)$ . First, assume that  $U'(0) \leq 0$ . Then, it follows from the concavity of the dictator's utility function that the  $\hat{a}_2$  maximizing the dictator's utility function on the interval  $[0, \rho\tau_D)$  is  $\hat{a}_2 = 0$ , which is independent of  $\rho$ .

Second, assume that  $U''(0) > 0$ , which happens for sufficiently high guilt aversion, i.e.,

$$\gamma_D > \frac{4}{12^k(\rho\tau_D)^{k-1}}.$$

Then, assuming a maximizer  $\hat{a}_3$  exists on  $[0, \rho\tau_D)$ , it must be an interior solution given by the following first-order condition:

$$\hat{a}_3 = \rho\tau_D - \left( \frac{4}{\gamma_D 12^k} \right)^{\frac{1}{k-1}}.$$

It can be seen that  $\hat{a}_3$  increases in  $\rho$ , in  $\tau_D$ , and  $\gamma_D$ , which proves **TP1**.

**Beliefs and Second-order Effects on Performance** So far we have implicitly assumed that the dictator's second-order beliefs are constant at  $\tau_D$ . However, as the dictator's equilibrium choice  $a^*$  weakly increases in  $\rho$ , first- and second-order beliefs must adjust accordingly. Hence,  $\tau_R$  and  $\tau_D$  must be weakly increasing in  $\rho$ . This yields **TP2**.

As the dictator's equilibrium action is increasing in  $\tau_D$ , this adjustment of beliefs leads to a (second-order) effect reinforcing **TP1**.

**Opt-out Decision** A (risk-neutral) recipient will opt in if

$$E[\pi_R] = 12E[\rho\tau_R(\rho)] > 2.$$

If  $\mu = 0$ , the dictator will choose  $a = 0$  and the recipient's beliefs will adjust accordingly. Hence,  $\tau_R(\cdot)$  and therefore also  $E[\pi_R] = E[\rho\tau_R(\rho)]$  will be higher for  $\mu = 1$  than for  $\mu = 0$ . As a result, opt-out rates are lower for  $\mu = 1$  than for  $\mu = 0$ . This yields **TP3**.

Example: Assume that  $\tau_R(\cdot) = 0$  for  $\mu = 0$ , and  $\tau_R(\cdot) = 1$  for  $\mu = 1$ . Then we would have  $0 > 2$  for  $\mu = 0$  and  $12\binom{1/2}{1/2}\binom{1/6}{5/6} = 6 > 2$  for  $\mu = 1$ .