

# Social Risk and the Dimensionality of Intentions

Jeffrey V. Butler<sup>a</sup>  
*EIEF and UNLV*

Joshua B. Miller<sup>b \*</sup>  
*Bocconi University and IGIER*

This version: *October 19, 2014*

## Abstract

Previous research has documented a behavioral distinction between “social risk” and financial risk. Individuals tend to demand a premium on the objective probability of a favorable outcome when that outcome is determined by a human counter-party instead of by a randomizing device (Bohnet, Greig, Herrmann, and Zeckhauser 2008; Bohnet and Zeckhauser 2004). Two explanations for this social risk premium have been offered: i) an aversion to a counter-party’s potentially malign *intentions*, or ii) a more general aversion to ceding control to another human with conflicting interests, irrespective of intentions (Bartling, Fehr, and Herz forthcoming; Bohnet and Zeckhauser 2004). An implication of the latter view is that social risk should always be aversive when the involved parties’ interest are at odds. In this paper we test for these two explanations experimentally by varying the degree to which outcomes reflect the counter-party’s intentions. Our study employs a between-subjects experimental design implementing slight modifications of the binary trust game of Bohnet and Zeckhauser across treatments. Our data support the first view, that intentions are a crucial determinant of the social risk premium. Intriguingly, we identify factors that can eliminate, or even change the sign of, the social risk premium. This result has implications for optimal contract design in a wide variety of situations involving social risk. In the penultimate section, we provide one explanation for this unexpected result which draws on the stereotype content model from the social psychology literature (Fiske, Cuddy, and Glick 2007).

**JEL Classification:** Z1, C91, D81

**Keywords:** Social Risk, Social Perception, Intention, Betrayal Aversion, Trust

---

\*a: Einaudi Institute for Economics and Finance, b: Department of Decision Sciences and IGIER, Bocconi University. Miller gratefully acknowledges support from a 2011 grant awarded by the Einaudi Institute for Economics and Finance. Both authors would like to thank seminar participants at EIEF and Bocconi University, as well as conference participants at MBEES in Maastricht (2014), FUR in Rotterdam (2014), and the ESA World Meetings in Honolulu (2014).

# 1 Introduction

The canonical framework for describing the domain of decision making under risk is that of a lottery, i.e., a probability distribution over consequences. However, a growing body of research investigating how risk and uncertainty affect behavior argues that an individual’s willingness to accept risk depends on factors other than merely probabilities and consequences. For example, several studies have noted that the source of risk, e.g., whether or not risk exposure is voluntary or not, affects decision-making and, at the same time, does not fit neatly within the consequentialist lottery framework (Loewenstein, Weber, Hsee, and Welch 2001; Slovic 1987).

In this paper, we focus on one source of risk that has recently captured economists’ attention: “social risk.” A decision maker faces social risk when another human being is the primary source of uncertainty (Bohnet et al. 2008). In their seminal contribution, Bohnet and Zeckhauser (2004), hereafter referred to as simply BZ, the authors use laboratory experiments controlling for many plausible extraneous factors—central among them, distributional preferences and ambiguity—to demonstrate that people treat social risk differently than an inanimate source of risk even when these two sources of risk implement identical probability distributions over monetary consequences. Specifically, in a situation where exposing oneself to social risk can result in either a monetary gain or a monetary loss, BZ document that individuals demand a premium in the probability of receiving the gain outcome in order to have uncertainty resolved by a human agent rather than by a randomizing device.<sup>1</sup> BZ attribute this *social risk premium* to *betrayal aversion* because it can be explained by an individual anticipating an additional disutility from an unfavorable outcome being chosen by a human agent, who can betray one’s “trust,” rather than by a randomizing device, which ostensibly cannot betray.<sup>2</sup> The betrayal aversion phenomenon has been documented in several subsequent studies conducted by a variety of authors across multiple cultures (Aimone and Houser 2012; Bohnet et al. 2008; Bohnet, Herrmann, and Zeckhauser 2010; Fetchenhauer and Dunning 2009, 2012).

The additional aversiveness social risk has been shown to engender may have a lot to do with the role of intentions. Obviously a human agent can have intentions, and can be perceived as acting intentionally, while a random device cannot. On the other hand, an alternative explanation, suggested by BZ themselves and seemingly supported by a growing number of closely related studies

---

<sup>1</sup>Bohnet and Zeckhauser (2004) elicit probabilities from first movers in a binarized version of the trust game experimental paradigm of Berg, Dickhaut, and McCabe (1995); the elicitation amounts to a quasi-strategy method applied to pooled second-mover responses. For earlier versions of incentivized experiments involving trust see Camerer and Weigelt (1988) and Fehr, Kirchsteiger, and Riedl (1993).

<sup>2</sup>On the surface, in BZ’s binary trust game the principal does not have sufficient information to know if the agent intends to betray, or if instead the agent simply has the selfish desire to obtain the higher payoff, viewing the consequences to the principal as simply an undesired (or neutral) side effect. In practice, people tend to attribute intentionality to harmful side-effects of selfish actions (Knobe 2006).

(Bartling et al. [forthcoming](#); Humphrey and Mondorf [2014](#); Neri and Rommeswinkel [2014](#); Owens, Grossman, and Fackler [forthcoming](#)), is that the premium has little to do with intentions *per se* but rather can be attributed to a more general aversion to ceding control to another human.

Separating these two explanations is made difficult by the fact that intention itself is a nuanced and multi-faceted concept (Bratman [1984](#); Mele [1992](#)). Intention need not translate into action and, moreover, the relationship between intention, action and outcome can be indirect. A person who intends to commit murder and takes all of the necessary actions (e.g., aims and pulls the trigger) but, by random chance, misses his target may be judged less harshly than somebody who did not intend to commit murder but, by random chance did (e.g., by driving drunk). At the same time, an amateur athlete may fully intend to score a goal or make a play, but because of inexperience or lack of ability the outcome may bear little resemblance to the athlete's intention. To make matters even more confusing, in the vernacular intention is often used as the opposite of action: I *intended* to donate blood, but never actually got around to it.

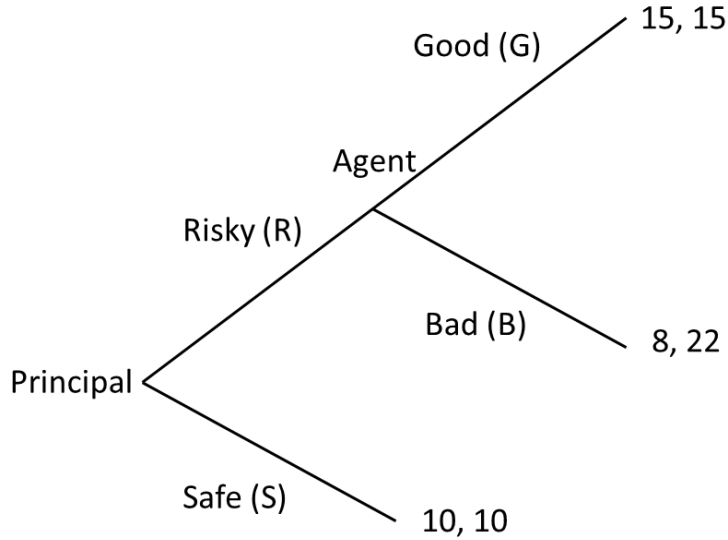
While the concept of intention does not have a universally agreed upon definition (Setiya [2014](#)), a handful of factors are commonly associated with the attribution of intentionality to actions and outcomes. For an individual to fully intend a consequence, the individual's actions must be voluntary, he or she must be able to foresee the consequences of the chosen action, and the individual must desire the consequence (Alicke [2000](#)). In line with this, to investigate directly the role intentions play in shaping the social risk premium documented by BZ and others we decompose intention into three components in our study: Act, Foresee and Desire. That is, an individual whose action determines the outcome, foresees the consequences of this action and, moreover, knows which outcome he or she desires can be said to act with full intention. In situations where the decision maker lacks one or more of these aspects, an individual's actions may be seen as reflecting the individual's intentions to a lesser extent, if at all.

Our approach will allow us to disentangle the impact of multiple facets of intention on the social risk premium in a way that cannot be done with the canonical betrayal aversion design since BZ did not construct their experiment with this goal in mind. In BZ's design, the ability to act and foresee are manipulated together: in their experimental treatment involving social risk, the agent acts, foresees and desires, while in their experimental treatment where an inanimate random device generates all risk, the agent lacks all but the last capacity. In our design, we employ the binary trust game and payoff parameters from BZ but experimentally manipulate multiple facets of an agent's ability to fully intentionally influence outcomes. This allows us to explore in more detail the role that intentions may play in determining attitudes toward the social dimension of risk than has previously been possible. Moreover, our design provides a straightforward test for two

fundamental open questions about social risk preferences. First of all, we can test whether social risk is *per se* aversive, or rather, like financial risk is sometimes even preferred (relative to baseline risk). Secondly, we will test whether the betrayal aversion phenomenon is merely a manifestation of a more general aversion to ceding control to another human.

Figure 1 displays the binary trust game from BZ which we also employ, with identical parameters and subject pools comparable to two of those found in Bohnet et al. (2008). One interpretation for this game is that it is a stylized representation of the interaction between a principal (first-mover) and an agent (second-mover). The principal decides whether to perform a task herself or to delegate the task to the agent. If the principal performs the task herself, she knows the payoffs that will result: 10 for both the principal and agent. Thus, not delegating is a S(afe) option from the principal’s perspective, involving neither social nor financial risk. If the principal delegates, the agent’s action can lead to one of two outcomes: earnings from the G(ood) outcome pareto dominate earnings from the non-delegation outcome, yielding earnings of 15 for both parties; earnings from the B(ad) outcome benefit the agent at the expense of the principal, yielding earnings of 22 for the agent but only 8 for the principal. Consequently, delegation contains an element of social risk, whose aversiveness may depend on intention. The possibility of disparate outcomes from delegation may stem from many sources. One plausible story is that the agent can exert effort which influences the outcome but is unobservable to the principal, rendering contracting infeasible (*cf.* Charness and Dufwenberg 2006; Charness and Levine 2007). We investigate how the agent’s capacity to intend the outcome, conditional on being allowed to determine the outcome, affects the principal’s decision to expose herself to social risk by varying (across treatments) the three fundamental facets of intention described above.

Consideration of these three facets naturally suggests four experimental treatments. We vary not only the agent’s ability to act, as in the original study, but also the agent’s ability to foresee and desire the consequences of his actions. We label the baseline binary trust game as Treatment AFD, because the agent can act voluntarily by choosing between G and B, can perfectly foresee the consequences of his action—(15,15) if he chooses G or (8,22) if he chooses B—and can desire the consequences. The three additional treatments involve removing one or more of the three facets of intention: the agent’s ability to act (A), to foresee the consequences of his actions (F), or his ability to desire the consequences of his actions (D). We label these treatments mnemonically as Treatment xxD, AxD, and Axx respectively, where “x” denotes the missing facet. In Treatment xxD, agents know all the outcomes possible in the game, and hence can have preferences over these outcomes (Desire) but, conditional on the principal choosing R, the outcome is completely determined by an inanimate randomizing device. In treatment Axx, the agent must choose between two options



**Figure 1:** *The Game*

without knowing anything about the consequences of his choices, not even that he is involved in a game. Finally, in Treatment AxD, the agent knows all of the possible outcomes in the game so that the agent can form preferences over outcomes and, moreover, the agent’s action determines which outcome obtains. However, we insert uncertainty between the agent’s action and the game’s outcome so that the agent cannot perfectly foresee the consequences of his actions.<sup>3</sup>

We implement a between-subjects design comprised of four separate treatments, with participants randomly assigned to treatments. Compared to BZ, our Treatment xxD is a version of their “Risky Dictator” game, while our AFD treatment is identical to their binary trust game. Our other two treatments (AxD and Axx) are novel and intended to test specific hypotheses. In each experimental session, each participant is randomly assigned to exactly one of these four treatments and is not informed of the existence of the other three treatments. All four treatments features a suitably modified version of BZ’s conditional binary trust game, which is one-shot, anonymous, and uses payoff parameters identical to the original study.

We randomly match principals with one agent drawn from a pool of potential agents. We use the strategy method for the agents, who choose between G and B before knowing if the principal with whom they have been matched has selected the safe option, S, or the risky option R. The principal’s decision between R and S is made conditionally. We ask the principal to state the minimum acceptable probability (MAP),  $p$ , of a randomly selected agent choosing G (or having G

<sup>3</sup>Technically, there are eight possible combinations of Act, Foresee and Desire. We did not find AFx or xxx of interest. We did not consider xFD as practically distinguishable from xxD. Finally, we do not anticipate xFx having any economically meaningful content.

selected for them) that would make the principal prefer R to S. If the relative frequency of agents choosing G (or having G selected for them) is greater than or equal to  $p$ , the principal commits to choosing R; otherwise the principal chooses S. This mechanism provides the principal proper incentives to truthfully report his or her MAP under mild assumptions which we discuss later.

Generally speaking, since a given MAP implements identical lotteries over outcomes across our four treatments, if social risk or intentions play no role in the principal’s preferences toward uncertain outcomes MAPs should not vary across treatments. If a simple aversion to ceding control to another human is the primary driver of differences in attitudes toward social risk, then since Treatments Axx, AxD, and AFD all involve ceding control to another human while Treatment xxD does not, we would expect MAPs to be similar across the former three treatments, all being substantially higher than in Treatment xxD. On the other hand, if an aversion to ceding control is not the entire story and intention is an important determinant of principals’ attitudes towards social risk, we would expect MAPs to vary over Treatments Axx, AxD, and AFD: while each of these treatments involves ceding control to a human agent, the agent’s capacity to intend the outcome of his action varies substantially across these treatments. We outline our specific hypotheses in a later section, after providing more detail on our experimental design.

As a preview of our findings, we replicate the original BZ results: on average, MAPs in Treatment AFD exceed those in our xxD treatment. At the same time, we find that MAPs vary across our treatments in a way that is not consistent with a simple aversion to ceding control. In particular, we identify a situation in which individuals actually have a relative *preference* for ceding control and exposing themselves to a social source of risk. In Treatment AxD, where we interfere with the agent’s ability to perfectly foresee the consequences of his actions, principals are on average willing accept a substantially lower probability of the Pareto-improving outcome G than in any of the other treatments. We provide an explanation consistent with this somewhat surprising finding in a later section.

Our study contributes to several literatures. The source of risk literature has investigated how an individual’s perception of and attitude towards risk is source dependent (Slovic 1987), and can be influenced by transient anticipatory emotions (Loewenstein et al. 2001; Slovic, Finucane, Peters, and MacGregor 2004). Recent work on the measurement of risk attitudes has found it useful to distinguish different sources of risk (Blais and Weber 2006; Dohmen, Falk, Huffman, and Sunde 2011; Weber, Blais, and Betz 2002). We contribute to this literature by showing that attitudes toward social sources of risk are context specific and more nuanced than previously thought.<sup>4</sup> In

---

<sup>4</sup>Note there is a literature on the social amplification of risk, which considers social processes rather than the mere fact that risk is generated by a human being (e.g. see Kasperson, Renn, Slovic, Brown, Emel, Goble, Kasperson, and Ratick (1988)).

particular, our findings suggest that individuals will not always demand a premium for exposing themselves to social risk, and that the scope for intentions, rather than the mere presence of a human decision maker with contrary preferences, is the likely determining factor. This implies that, for example, adding social risk on top of financial risk to an investment decision by delegating investment authority to a financial intermediary with potentially conflicting interests may not always reduce an individual’s propensity to invest even when the intermediary contributes little in the way of expertise or knowledge.

We also contribute to the literature on intentions. While the seminal works in this vein (Dufwenberg and Kirchsteiger 2004; Rabin 1993) predict that intentions should matter, this prediction has typically been tested using designs that take away intentions by removing an individual’s ability to act altogether, often by substituting a random device (Blount 1995; Charness 2004; Falk, Fehr, and Fischbacher 2008) or a neutral third party (Charness 2000; Kagel and Wolfe 2001). The betrayal aversion literature follows this schema, typically removing intention in some treatments by having a random device determine outcomes. In such designs, not only can an agent no longer intend a specific consequence, he cannot even influence the consequence, as he is a passive stakeholder rather than a counter-party. One interpretation of these results is that, rather than caring about intentions *per se*, principals are sensitive to consequences in social games where an agent has an alternative course of action that would lead to a different payoff (as in Gurdal, Miller, and Rustichini (2013)). The contrast between our Treatments AFD and xxD, which is essentially identical to the comparison made in BZ, has this very issue. The principal may be anticipating intentional betrayal, or the principal may be anticipating a counterfactual comparison. Unlike Treatment xxD, Treatment AxD allows the agent to be active in determining consequences by providing the agent with alternative courses of action while at the same mitigating intentional “betrayal:” the agent cannot fully intend a consequence he cannot foresee. Our finding that principal’s MAPs are lower in AxD than in AFD suggests that intentions are an important determinant of the social risk premium and hence the betrayal aversion phenomenon. Furthermore, comparing AxD with xxD allows us to check if individuals treat a situation where a human agent cannot fully intentionally betray them (AxD) as identical to a payoff-equivalent situation where a random device chooses on behalf of the agent (xxD). We find that principal’s MAPs are significantly lower in AxD than in xxD. This suggests that individuals may be willing to accept a discount for exposing themselves to a social risk when the risk is generated by someone whose ability to intentionally harm them is constrained but not fully eliminated. Interestingly, this somewhat puzzling finding is consistent with recent results in the social psychological literature, which we discuss in our concluding remarks.

The most closely related study to ours in the literature focusing on intentions, and to the best

of our knowledge the lone exception to the criticism above, is Charness and Levine (2007). There, the authors study the interaction between a principal (firm) and an agent (worker). The firm moves first by setting a wage, which is then either increased or decreased by random chance. The worker observes the firm’s wage choice and the realization of the chance move, and then decides whether to punish or reward the firm with (costly) high or low effort, or to exert medium effort at no cost to either the firm or the worker. The data suggest that worker effort responds to the valence of a firm’s intentions (good or bad) independent of the consequences they generate. This provides some support for the idea that there is some scope for intentions to matter in our AxD treatment, where a chance move is inserted between an agent’s choice and the outcome that is implemented. In contrast to Charness and Levine, however, our study focuses on how the agent’s capacity to intend affects the principal’s willingness to expose herself to the social risk, when the valence of intentions is unknown.

We contribute, as well, to a growing literature on the intrinsic value for control. An early conjecture made, but not directly tested, by Bohnet and Zeckhauser themselves (2004, p. 478) is that the social risk premium may simply be one manifestation of a more general aversion to social sources of risk that is driven by a basic desire to avoid relinquishing control to another human. Recent research in this vein has in fact been pointing toward a growing consensus on just this point: that ceding control is generally aversive. Bartling et al. (forthcoming), for example, vary the conflict of interest between two parties and explicitly elicit individuals’ valuations for retaining decision-making authority, i.e., control. They find that this value is positive and significant on average. In a related study, Neri and Rommeswinkel (2014) also estimate a positive value for retaining control. Using a design similar to ours, Humphrey and Mondorf (2014) argue that ceding control to an agent lacking the ability to intend is as aversive as ceding control to an agent who can fully intend his actions, suggesting that the betrayal aversion phenomenon is primarily about an aversion to ceding control.<sup>5</sup> Owens et al. (forthcoming) estimate a substantial positive value for control over own earnings even in the absence of conflict of interest. In our study, we keep the severity of conflict of interest fixed and vary the ability for an agent to intend outcomes determined by his actions. In contrast with the existing studies, we find that ceding control is not always

---

<sup>5</sup>Their design has important differences from ours. First, they purport to take away an agent’s ability to intend by implementing G or B based on whether the agent can correctly guess whether the principal’s birth year was even or odd. In this setting, the principal could plausibly believe that the agent could guess odd or even with more than 50 percent accuracy: in any given year, or any given university cohort, odd and even birth years are not uniformly distributed. This leaves some scope for the principal to believe that the agent can intentionally implement an outcome, and these beliefs may vary across principals. Consequently, it is not clear whether their null finding stems from the fact that intentions do not matter or from a failure to eliminate the agent’s ability to intend. Secondly, different treatments were conducted on different days rather than randomizing into treatments within each session, so that selection effects or day-specific fixed effects may introduce unwarranted noise which could contribute to their null finding.



aversive even when preferences are misaligned.

The remainder of the paper proceeds as follows. First, we describe our experimental design and procedures in detail. Next, we specify our hypotheses. Subsequently, we present our results. In the penultimate section we describe one potential unifying explanation for our results drawing on recent research from the social psychological literature. In the concluding section we discuss our findings, putting them in context and highlighting avenues for future research.

## 2 Experimental Design and Procedures

All experimental sessions were conducted at Bocconi University. The participants consisted of 158 undergraduate students recruited from the Bocconi Experimental Laboratory for the Social Sciences (BELSS) on-line subject recruitment system and 11 graduate students who were recruited individually via email.<sup>6</sup> Neutral wording for roles and outcomes were used in the experimental instructions (see Appendix B). Here, we use more descriptive terminology for ease of exposition.

All undergraduate students were assigned the role of principal (described below). Each undergraduate student participated in a single session of approximately 27 students in “Room X” (6 sessions total) and was randomized into exactly one of four treatments, again described below. Each graduate student in “Room Y” was assigned the role of agent and was assigned to all four treatments of each “Room X” session.<sup>7</sup> Our focus will be on participants assigned the role of principal, all of whom were informed only that they were matched with participants in Room Y, on the fifth floor of the building, whose decisions could potentially affect their own outcomes.

Upon arrival in Room X, principals selected one card from a shuffled deck of cards. The chosen card revealed a “code” number, from 1 to 27. Participants were instructed to sit in the private carrel of the laboratory corresponding to their code numbers. Once all participants were seated they were immediately informed that the amount they would be paid at the end of the session would be between 8 Euros and 15 Euros. Furthermore, they were informed that how much they would be paid depended on a single choice that they would make together with a choice of the agent in Room Y with whom they had already been matched, and that they should therefore listen and read carefully.<sup>8</sup>

Next the experimenter read a welcome script (See Appendix B.2). Principals were informed that they were in Room X and the code number they drew upon entering the room determined the

---

<sup>6</sup>The website and database is administered by Sona-Systems. The average age of those in our sample was 21 students.

<sup>7</sup>The agents submitted decisions in each treatment in the following order: Axx, AxD, AFD, and xFD. Agents received payments from all four treatments, but principals were not informed of this.

<sup>8</sup>The participants in Room Y decided in advance using the strategy method as is detailed in the treatment description below.

agent they were matched with in Room Y. They were told that all participants’ identities would be kept anonymous.<sup>9</sup> Next, principals were given an overview of the experimental session and told that: (1) they would receive detailed instructions for them to read to themselves; (2) they would answer a short quiz intended to check their understanding of the instructions; (3) we would check their responses to the quiz; (4) they would respond to the single “Key Question” (which was our outcome measure of interest); (5) while we were matching their responses to that of the other player in the role of agent from Room Y, they would fill out a survey; and finally, (6) that they would be paid.

For each X-session block, the shuffled “code” numbers implemented a permuted block randomization with a uniform allocation of participants into treatments.<sup>10</sup> Principals were not aware they were assigned to a treatment nor that there were other treatments.<sup>11</sup>

**The Treatments** Each treatment had a common underlying structure in terms of how payoffs depended on the choices of the principal and the agent (Figure 1).<sup>12</sup> In all treatments, the principal made a decision which determined whether S, the safe option, or R, the risky option, was selected. If the principal “chose” option S, this yielded a certain outcome for the pair, and the experimental earnings were 10 euros to the respective players. If the principal “chose” option R, then the payoffs depended on the choice of the agent. If the agent chose option G, then experimental earnings were 15 euros for both players. If the agent chose option B, then the agent received 22 euros in experimental earnings and the principal received 8 euros in experimental earning.<sup>13</sup> In all treatments the agent’s choice between options G and option B conditional on the principal choosing alternative R was determined before the principal chose (strategy method). In each treatment, the principal’s choice between alternatives S and R was implemented as in the trust game of BZ with instructions that closely followed those reported in Bohnet et al. (2008). For each subject in the role of principal, we elicited the value  $p$  which was described as their *minimum acceptable probability* (MAP) of being matched with an agent choosing option G that would lead them to choose alternative R over alternative S. If the percentage of subjects in Room Y choosing G,  $p^*$ , was greater than or equal to  $p$ , then the principal committed to choosing alternative R and payoffs would be determined by

---

<sup>9</sup>It was important to match before they choose rather than after, to make clear the decision is coming from their assigned agent rather than the outcome being a draw from a pool of already determined decisions.

<sup>10</sup>For sessions of 27 participants the allocations in each treatment were 6,7,7,7. For sessions of 26 participants the allocations in each treatment were 6,6,7,7.

<sup>11</sup>The agents were provided identical instructions as the principals, and given a separate sheet to submit their responses. In Treatment Axx the agent had no instructions and only a choice

<sup>12</sup>Payoffs were denominated in euros, the numerical amounts were identical to the original experiment Bohnet and Zeckhauser (2004) as well as the as well as the cross-cultural studies Bohnet et al. (2008) and Bohnet et al. (2010).

<sup>13</sup>In the experiment, we use the more neutral letters J and K for options G and B, respectively, primarily because they are not present in the Italian alphabet and thus participants were unlikely to have associations with these letters.

the choice of the agent with whom the principal had (already) been matched.<sup>14</sup> If we assume that the difference in utility between S and R as a function of  $p^*$ ,  $\Delta U(p^*) = U(R, p^*) - U(S, p^*)$ , is strictly increasing in  $p^*$  on  $[0, 1]$  and satisfies  $\Delta U(1) > 0 > \Delta U(0)$ , then for the principal reporting  $p = MAP$ , where  $MAP := \inf\{p^* > 0 : \Delta U(p^*) > 0\}$ , is a unique weakly dominant strategy.<sup>15</sup>

The random matching made each treatment structurally identical for a principal who cares only about the probability that the agent “chooses” G. The distinguishing feature between treatments was the mode in which the agent “chooses” between options G and B. These differences are presented below for each treatment:

*AFD*: Each agent decided directly between option G and option B. This is the Trust Game from Bohnet and Zeckhauser (2004).<sup>16</sup>

*Axx*: Each agent was presented with a row of 17 cells and asked to choose one cell. Agents were not aware that they were playing a game with another player, and therefore they did not know that in each cell there was either a G or a B, and that the cell they selected would determine the payoff for themselves and for their respective principals. The principals were aware of this.<sup>17</sup>

*AxD*: Each agent was presented with a row of 17 cells and asked to choose one cell. Agents were aware that in each cell there was either a G or a B, and that the cell they selected would determine the payoff for themselves and for their principals. The agent could not see the contents of each cell and thus could not foresee whether G or B would be selected. Principals also could not see the contents of the cells. Both principals and agents were aware of this.

*xxD*: Each agent was presented with a row of 17 cells and a randomizing device determined one of the cells at random.<sup>18</sup> Agents were aware that in each cell there was either a G or a B,

---

<sup>14</sup>By stating the question in this way, our design, and that of Bohnet and Zeckhauser (2004), implicitly assume that principals would like to statistically discriminate against agents. If principals did not wish to discriminate they could simply report  $p = 0$  or  $p = 1$ , and 4 out of 158 subjects did this.

<sup>15</sup>Alternatively, using the justification present in Bohnet et al. (2008), if we assume players believe that  $p^*$  is drawn from a distribution where the support contains a neighborhood of their MAP, then reporting  $p$  equal to their MAP is strictly dominant as this is equivalent to a Becker-DeGroot-Marshak (BDM) elicitation procedure with  $p^*$  generated by the collective behavior of agents in Room Y (Becker, Degroot, and Marschak 1964).

<sup>16</sup>The instructions for Treatment AFD can be found in Appendix Section B.3.

<sup>17</sup>The instructions for the principals were identical to the instructions in treatment AFD, except for the existence of the 17 cells. The choice of 17 cells was made so that a uniform distribution over G and B was unlikely to be focal from the perspective of the principals. The decision not to include the 17 cells in treatment H when implementing the design was made to avoid the risk that subjects could become confused with the introduction of a device that transparently serves no purpose. While the visual representation of the instructions is slightly different, there is some indication that it is justifiable to assume that these differences are slight, and do not confound our results: the average MAP in treatment xxD is only slightly higher than that in the study of Bohnet et al. (2010), which has a comparable pool of subjects (see Table 2).

<sup>18</sup>The selection was made using [www.random.org](http://www.random.org).

and that the cell the randomizing device selected would determine the payoff for themselves and for their principals. The agent could see the contents of the cells, but the principal could not. Both principals and agents were aware of this.<sup>19</sup>

### 3 Hypotheses

Denote by  $MAP_{AFD}$  the principals' average MAP in Treatment AFD, and the other treatments' average MAPs analogously. In line with previous research suggesting an important role for intentions in situations like the one studied here where one player can help or harm another player, our design generates four natural hypotheses.

*Hypothesis 1:  $MAP_{AFD} > MAP_{xxD}$ .*

In Treatment xxD, the agent is entirely passive. A randomizing device with a fixed probability selects between  $G$  and  $B$  on behalf of each agent before the principal chooses. In Treatment AFD the outcome fully reflects the agent's intentions. Treatments xxD and AFD together essentially replicate the canonical betrayal aversion setup of BZ, although Treatment xxD is slightly modified to keep these two treatments more parallel than in BZ.<sup>20</sup> Hypothesis 1 states that we expect to replicate the phenomenon which has been labeled as "betrayal aversion" in our modified setting. One explanation for this phenomenon is that principals anticipate an additional disutility from the bad outcome B in AFD relative to xxD because in the former this outcome fully reflects the agent's intention.

*Hypothesis 2:  $MAP_{AxD} < MAP_{AFD}$ .*

Hypothesis 2 states that the agent's ability to foresee the consequences of his actions is important in determining the aversiveness of intentions. In utility terms, this would be consistent with the principal anticipating that B will *feel worse* in AFD than in AxD, and hence yield lower utility in AFD than in AxD, because in the former the harm to the principal is fully intended while in the latter the harm may reflect intentions only partially. It would also be consistent with the principal anticipating that G will *feel better* in AxD than in AFD, yielding higher utility in the former than

---

<sup>19</sup>Unlike the Risky Dictator game reported in BZ, which treated  $p^*$  as an ex-ante probability of a random device, in all treatments of our study  $p^*$  was equal to the empirical relative frequency of G choices, even if the G choices were determined by realizations of the randomizing device.

<sup>20</sup>In our setup, in both xxD and AFD principal's MAPs are compared to the realized relative frequency of  $G$  in the population of agents. In contrast, in BZ's risky dictator game MAPs were compared to the *ex-ante* probability of their randomizing device selecting the outcome B, while in their binary trust game MAPs were compared to the realized relative frequency.

in the latter. In the two-outcome setting we study these effects are not separately identifiable. However, we provide a rationale for the latter effect in a later section.

*Hypothesis 3:  $MAP_{xxD} = MAP_{Axx}$ .*

In Treatment Axx the agent’s action determines the outcome, but the agent knows nothing about the (strategic) situation he is in. The agent is essentially a human randomizing device. Hypothesis 3 therefore states that the mere fact that control over outcomes is ceded to a human agent instead of an inanimate random device will not influence how aversive the principal finds the situation.

*Hypothesis 4:  $MAP_{AxD} > MAP_{xxD}$ .*

The primary difference between Treatment AxD and Treatment xxD is whether the agent can take an action which influences the outcome. Therefore, Hypotheses 4 states that ceding control to a human agent with conflicting interests, even when that agent cannot perfectly foresee the consequences of his actions (AxD), is more aversive than a situation where the agent cannot influence the outcome at all but still knows which outcome is personally desirable (xxD).

## 4 Results

In Table 1 it is evident that the average MAP in our Treatment AFD is similar to those reported in similar studies for other western countries (Switzerland and the United States) in the most directly related treatment of Bohnet et al. (2010).

**Table 1:** *Minimum Acceptable Probabilities in Treatment AFD (Mean, Median, [N])*

	ALL	Men	Women
Milan	0.54 <i>0.55</i> [38]	0.51 <i>0.45</i> [25]	0.59 <i>0.60</i> [13]
Switzerland	0.51 <i>0.55</i> [25]	0.46 <i>0.48</i> [18]	0.62 <i>0.60</i> [7]
United States	0.54 <i>0.50</i> [31]	0.50 <i>0.50</i> [19]	0.61 <i>0.72</i> [12]

Data from Switzerland and the United States are from the Trust Game of Bohnet et al. (2010).

Our xxD treatment employs a randomizing device, rather than the Room Y co-player, to determine the selection between G and B. This treatment is analogous to the risky dictator game of

Bohnet et al. (2010).<sup>21</sup> In Table 2 we see that the MAPs in our study are not as low as in Bohnet et al. (2010). A potential explanation for this is that the risk generated from being randomly matched to an outcome from a sample of realizations of a random device is perceived differently than the equivalent risk of receiving a single outcome directly from the device itself.<sup>22</sup> The results in Table 2 suggest that the Betrayal Aversion effect itself could be (partially) driven by the effect of these differences.<sup>23</sup>

**Table 2:** *Minimum Acceptable Probabilities in Treatment xxD (Mean, Median, [N])*

	ALL	Men	Women
Milan	0.46	0.43	0.54
	<i>0.50</i>	<i>0.40</i>	<i>0.60</i>
	[40]	[27]	[13]
Switzerland	0.40	0.33	0.48
	<i>0.42</i>	<i>0.30</i>	<i>0.50</i>
	[24]	[13]	[11]
United States	0.32	0.28	0.38
	<i>0.29</i>	<i>0.29</i>	<i>0.35</i>
	[29]	[16]	[13]

Data from Switzerland and the United States are from the “Risky Dictator Game” of Bohnet et al. (2010).

**Comparisons Across Treatments** Having established broad comparability with previous results we now turn to comparing patterns across treatments within our own study. In Table 3 we present simple means of participants’ MAPs for each our treatments, separately. The main statistical test we employ is the permutation test, as random assignment into treatments was conducted as a permuted block design, with stratification at the session level. We report means and  $p$ -values resulting from each of the six possible pair-wise tests in Table 4

Comparing the average MAP in our AFD treatment to MAPs in our xxD treatment (Table 3), we find evidence consistent with previous results on betrayal aversion.

<sup>21</sup>In our study, however, MAPs are elicited in a more parallel fashion across treatments: MAPs are compared to the actual empirical relative frequency of G choices ( $p^*$ ) made by a randomizing device on behalf of the agents, rather than to the ex-ante probability of G being selected by the randomizing device itself.

<sup>22</sup>In Bohnet et al. (2010) the device was a an urn with an unknown distribution of balls, in our study it is the third party uniform random number generator based on atmospheric noise ([www.random.org](http://www.random.org)) and an unknown distribution of Gs and Bs to select from. Some potential sources of the perceptual difference could be: (1) whether chance is realized ex-ante to the decision or ex-post (citations needed); (2) the mechanism itself (matching vs. draws), which could be a real issue because that means that the measurement device creates an artifact (our evidence suggests it is not an artifact, which is a potentially important ancillary result); (3) vulnerability to experimenter manipulation.

<sup>23</sup>Or if one is unwilling to view our 1-17 representation as innocuous, then the difference could be there.

*Result 1: Hypothesis 1 finds support in our data. Principals' MAPs in our AFD treatment are larger on average than MAPs in Treatment  $xxD$*

In our study, the raw averages in (Table 3) suggest that principals are willing to pay a 7 percentage point MAP premium to have outcomes determined by an inanimate randomizing device instead of a human agent who has conflicting monetary interests and can perfectly implement his desired outcome. This result provides additional evidence for the robustness of the original BZ findings and, at the same time, provides reassurance that our experimental design and subject pool are reasonable. The difference calculated is statistically significant ( $p = 0.048$ ) using a formal non-parametric one-sided permutation test which matches our experimental randomization procedure (stratified by session).

**Table 3:** *Minimum Acceptable Probabilities Across Treatments (Mean, StdDev, [N])*

	AFD	AxD	Axx	xxD
MAP	0.54	0.37	0.50	0.47
	<i>0.03</i>	<i>0.04</i>	<i>0.04</i>	<i>0.03</i>
	[38]	[40]	[38]	[40]

Next, consider how principals' attitudes toward the social risk embodied in Treatment AFD compares to principals' attitudes toward the social risk stemming from our AxD treatment. Recall that in both of these treatments a human agent to whom decision authority has been ceded would take an action which decides the outcome. The primary difference between these treatments is whether this agent can perfectly foresee the consequences of his action. If the ability to foresee consequences is an important component of intention, and intention matters for social risk preferences, we should observe lower MAPs in Treatment AxD than in Treatment AFD.

*Result 2: Hypothesis 2 finds support in our data. Principals' MAPs in our AxD treatment are smaller on average than MAPs in our AFD treatment.*

On average, MAPs are 14 percentage points lower in AxD than in AFD (Table 3). This probability premium is substantial, being twice as large as the probability premium associated with an inanimate randomizing device. The test reported in Table 4 reveals the premium is also highly statistically significant ( $p < 0.01$ ).

Our next comparison addresses directly whether the mere fact of having a human agent decide (Treatment Axx) as opposed to an inanimate randomizing device (Treatment xxD) leads to a difference in the principal's willingness to expose herself to risk.

*Result 3: Hypothesis 3 finds support in our data. Principals' MAPs in our Axx treatment do not differ significantly from MAPs in xxD.*

Referring to Table 3 once again, we find a small three percentage point difference in average MAPs between these two treatments. Table 4 suggests this difference is not statistically significant ( $p = 0.7881$ ). In addition to providing evidence about our third hypothesis, this non-significant difference provides reassurance that extraneous differences between human decision-makers and randomizing devices, such as the degree of ambiguity involved in their respective decision procedures, is not driving our results or previous results in the betrayal aversion literature.

Our final hypothesis rests on the notion that the capacity for counter-parties to have malign intent is always aversive, which would be in line with, but not necessarily implied by, existing research. In particular, since all existing studies we are aware of feature a comparison between treatments where outcomes fully reflect agents' intentions versus treatments where outcomes cannot reflect agents' intentions at all, the literature is mostly silent on how social risk attitudes respond to situations where outcomes plausibly reflect some part, but not all, of agents' intentions.

*Result 4: Hypothesis 4 does not find support in our data. Principals' MAPs in Treatment AxD are significantly lower than MAPs in our xxD treatment.*

In words, Result 4 indicates that principals are less averse to risk when it comes from a social source if the agent cannot foresee the consequences of his actions (the scope for intentions are limited) than when it comes from a non-social source like a random device. In this case, Principals are willing to accept a substantial 10 percentage point *lower* probability of the good outcome when ceding decision authority to the human agent relative to ceding decision authority to an inanimate randomizing device ( $p = 0.033$ , one-tailed permutation test).

To provide some evidence on whether this last surprising finding is a statistical fluke from conducting multiple hypothesis tests, notice that although we had no direct *a priori* hypothesis on the comparison between Treatments AxD and Axx, since we did hypothesize that the agent in Treatment Axx would be essentially equivalent to the random device in Treatment xxD, one might conjecture that the patterns in MAPs when comparing Treatments AxD to xxD would be similar to the patterns in MAPs from the AxD vs. Axx comparison. In fact, this is exactly what we find. The probability premium in MAPs associated with the comparison between these latter two treatments (AxD vs. Axx) has the same sign and is of similar magnitude to the premium associated with the former two treatments (AxD vs. xxD) and, moreover, is also highly statistically significant (13 percentage points;  $p < 0.01$ ). This provides some reassurance that Result 4 is a robust finding



and bolsters the interpretation that intention in general, and the capacity of the agent to desire (contrary) outcomes in particular, is an important determinant of attitudes toward social risk.<sup>24</sup>

**Table 4:** *Pair-wise between-treatment differences in the mean MAPs and p-values for permutation tests (one-tailed). The permutations test is stratified at the session level to match the permuted block design (See Appendix Section A for details).*

Comparison	Difference	P-Value
AFD vs. xxD	0.08**	.0483
AFD vs. AxD	0.17***	.0002
AFD vs. Axx	0.04	.1883
xxD vs. Axx	-0.04	.7881
xxD vs. AxD	0.09**	.0325
Axx vs. AxD	0.13***	.0038

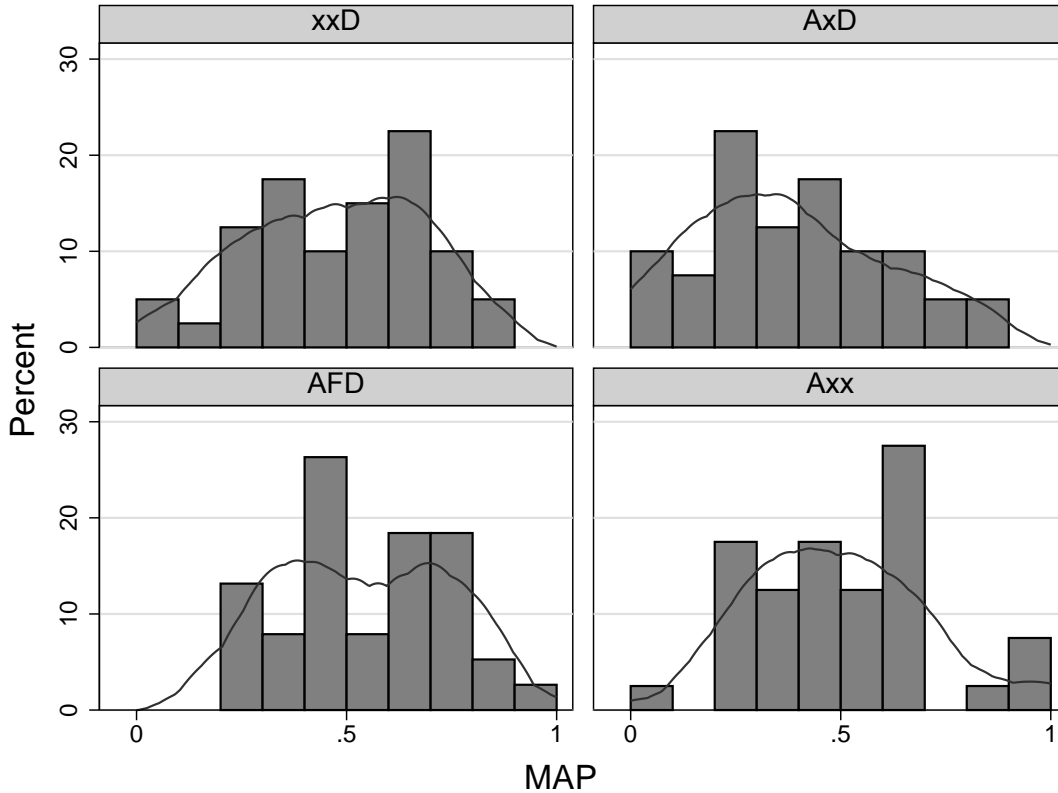
\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$  (one-sided, right)

Turning from simple means to distributions, in Figure 2 we present histograms, overlaid with kernel density estimates, of principals’ MAPs for each of our treatments separately. An important point to notice is that in all treatments a wide range of MAPs are reported. This is important because it makes it unlikely that our results are driven by a few outliers. A second point to notice is that the histograms and kernel density estimates tend to corroborate the story gleaned from comparing means. For example, the distribution of MAPs in our AxD treatment are essentially a leftward shift of the MAPs in our AFD treatment. Low MAPs (more trusting behavior) are more prevalent when we introduce noise into the mapping between co-player’s action and outcomes, bringing down the average MAP for the AxD treatment.

## 5 One possible explanation: competence

A theoretical framework exists in the social psychological literature that can explain both our surprising finding (Hypothesis 4) and the findings we correctly anticipated (e.g., Hypothesis 1). A recent theory of how people perceive strangers and form stereotypes about social out-groups, the *stereotype content model* (Fiske, Cuddy, Glick, and Xu 2002), incorporates insights from earlier work on social perception (Asch 1946; Bales 1950; Rosenberg, Nelson, and Vivekananthan 1968) and organizes patterns of how stereotypes operate in different cultures (Cuddy et al. 2009; Fiske et al. 2002). The basic finding is that the personality impressions that people form about other individuals, how they construe behavior, and the stereotypes they hold about members of other

<sup>24</sup>At this point it is worth emphasizing that while we report all pair-wise treatment comparisons for completeness, some of these pair-wise differences are difficult to interpret conceptually or theoretically because two conditions change at once. The findings in the following comparisons should therefore be interpreted cautiously: AFD vs. xxD; AFD vs. Axx; xxD vs. Axx.



**Figure 2:** Histograms of MAPs by treatment (with Kernel Density)

groups, can be categorized into two factor dimensions, which have been labeled *warmth* and *competence* respectively (Fiske et al. 2007). Here, warmth is largely synonymous with intent—positive (negative) intentions being identified with high (low) warmth.<sup>25</sup> The impression of another agent’s personal warmth serves as a cue to what the other agent’s possible goals are with respect to the self, while the impression of the agent’s competence is thought to serve as a cue to the agent’s ability to carry out those goals.<sup>26</sup>

Our experimental treatments can be re-cast in terms of these two dimensions: warmth and competence. For example, in Treatments AFD and AxD, conditional on being given decision authority the agent’s preferences are (equally) in conflict with the principal’s so that the agent’s “warmth” across these treatments should be similar. However, the agent’s ability to pursue his

<sup>25</sup>In their paper introducing the stereotype content model, Fiske et al. (2002) present a model of people as having pragmatic/consequentialist objectives when dealing with strangers: “when people meet others as individuals or group members, they want to know what the other’s goals will be vis à vis the self or in-group and how effectively the other will pursue those goals. That is, perceivers want to know the other’s intent (positive or negative) and capability; these characteristics correspond to perceptions of warmth and competence respectively.”

<sup>26</sup>In the fields of management and sociology there is a similar two-factor definition of trust: trusting another human agent involves (1) trusting in an agent’s competence, and/or, (2) trusting in an agent’s intentions (Nooteboom 2002).

preferences (*competence*) varies substantially across these treatments. In Treatment AFD the agent is fully competent while in Treatment AxD the agent's competence is quite limited. By way of contrast, in Treatment xxD since the agent takes no action person perception may not be implicated at all, these two factor dimensions being essentially irrelevant.

Assessing the risk generated by interacting with another human agent in an experimental social dilemma, such as the trust game, also involves an act of social perception in order to anticipate the behavior of one's counterpart. While an anonymous laboratory setting appears to provide little scope for forming personality impressions or using stereotypes, evidence suggests that the perceptions and stereotypes people form about others can be driven solely by context, in particular the degree of competition and the relative control over resources one's counterpart has vis-à-vis the self (Cikara and Fiske 2013; Fiske et al. 2002). Context can determine the formation of these personality impressions, which can, in turn, predict affective reactions. Social out-groups who compete for resources and successfully control them in their own favor tend to be viewed as having low warmth and high competence, which, in turn leads them to be envied or perceived as a threat. The prospect of ceding control to such agents may therefore generate a negative affective reaction (Fiske et al. 2002). On the other hand, if these same agents experience personal misfortune, a positive affective reaction may result if it is perceived as a removal of a social threat; this reaction has been associated with an emotion known as *schadenfreude* (Cikara and Fiske 2013). Since the probability of the agent's misfortune increases in his incompetence when that agent has control, the knowledge that the agent, competing for the same resources, may not be able to effectively pursue his own interests could counterbalance or even overcome the negative affect associated with the prospect of ceding control to him.

One explanation for our findings is that principals involved in social dilemmas respond to changes in these contextual details in a pattern consistent with the stereotype content model. A principal may exhibit an aversion to the possibility of betrayal and demand a premium to expose herself to social risk because of the negative *affect* associated with facing the potentially threatening intentions of a human agent who competes for resources and can competently control the outcome in his favor. By contrast, if a human agent can control resources but cannot competently do so in his own favor, then the principal is partially protected from the influence of any threatening intentions the agent may have. The principal may even delight in anticipating the possibility of the agent's malign intentions translating into (unintentional) personal misfortune for the agent. These considerations may, in turn, generate a positive affect towards the prospect of being exposed to social risk, and partially offset or even change the sign of the social risk premium.

## 6 Concluding Remarks

In this study we experimentally investigated the role of intentions in determining attitudes toward social risk. Social risk is an important and ubiquitous phenomenon, being present whenever decision authority is delegated from a manager to an employee or from an investor to a financial intermediary or advisor. Existing literature seemed to suggest that there would be little to learn from this investigation, that social risk should always add to the aversiveness of the underlying monetary risk involved and that only the underlying mechanism for this additional aversiveness was in question. We found that counter-party intentions may play a large role in determining attitudes toward social risk and that, contrary to all existing studies we are aware of, in some situations the presence of social risk may actually be preferable to outcome-equivalent purely financial risk. Our findings have broad implications for contract design in situations where the presence of social risk is a choice variable. One can imagine firms may choose to make financial advice automated (low social risk) or to deeply involve financial advisors with potentially conflicting goals (high social risk).

We went on to provide one explanation that is consistent with findings, drawing on recent research in social psychology. In situations where a counter-party may take decisions which can help or harm an individual, impressions about whether the counter-party has conflicting interests and about the counter-party's competence in pursuing his or her goals are subconsciously formed and produce affective responses guiding decision-making (Fiske et al. 2007). When interactions are anonymous, context alone may generate similar affective responses (Fiske et al. 2002). With this in mind, one interpretation of our treatments is that we held monetary conflict-of-interest constant across most of our treatments while varying the agent's competence. Treatment AFD featured high agent competence while AxD reflected low agent competence. Since an agent took no action in xxD such affective responses were unlikely to have been implicated. The prospect of ceding control to a fully competent agent with conflicting interests may have generated a negative affective reaction. On the other hand, if these same agents experience personal misfortune, a positive affective reaction may result if it is perceived as a removal of a social threat; this reaction has been associated with an emotion known as *schadenfreude* (Cikara and Fiske 2013). Since the probability of such misfortune increases in incompetence when an agent has control, the knowledge that an agent with conflicting interests may not be able to effectively pursue these interests could counterbalance or even overcome the negative affect associated with the prospect of ceding control to him.

Our results suggest avenues for future research. In the present experimental study, we externally manipulated the agent's competence across treatments. More generally, competence may be something that varies across agents and in which an agent may invest through costly skills acquisition. It is unclear what may happen when incompetence is an endogenous factor. Even maintaining

exogenous competence, however, a second question for future research suggests itself. Suppose the agent's competence is only imperfectly observed by the principal. One potentially interesting implication of our finding that agent incompetence may increase the principal's tolerance for social risk is that, knowing this, a strategic agent may feign incompetence to elicit delegation. Such "strategic incompetence" requires that agents correctly anticipate the effect of incompetence on principals' attitudes toward social risk that we document so that it is not clear how large a role feigned incompetence may play in actual behavior (*cf.* the "Lure" treatment of Charness, Rustichini, and Van de Ven 2013). Again, we leave for future research the introduction of asymmetric information about the agent's competence into the situation studied here.

## References

- AIMONE, J. A. AND D. HOUSER (2012): “What you don’t know won’t hurt you: a laboratory analysis of betrayal aversion,” *Experimental Economics*, 15, 571–588.
- ALICKE, M. D. (2000): “Culpable control and the psychology of blame,” *Psychological Bulletin*, 126, 556–74.
- ASCH, S. E. (1946): “Forming impressions of personality,” *Journal of Abnormal and Social Psychology*, 42, 248–290.
- BALES, R. (1950): “A set of categories for the analysis of small group interaction,” *American Sociological Review*, 15, 257–263.
- BARTLING, B., E. FEHR, AND H. HERZ (forthcoming): “The Intrinsic Value of Decision Rights,” *Econometrica*.
- BECKER, G. M., M. H. DEGROOT, AND J. MARSCHAK (1964): “Measuring utility by a single-response sequential method,” *Behavioral Science*, 9, 226–232.
- BERG, J., J. DICKHAUT, AND K. MCCABE (1995): “Trust, Reciprocity, and Social History,” *Games and Economic Behavior*, 10, 122–142.
- BLAIS, A.-R. AND E. U. WEBER (2006): “A domain-specific risk-taking (DOSPERT) scale for adult populations,” *Judgement and Decision Making*, 1, 33–47.
- BLOUNT, S. (1995): “When social outcomes aren’t fair: The effect of causal attributions on preferences,” *Organizational Behavior and Human Decision Processes*, 63, 131–144.
- BOHNET, I., F. GREIG, B. HERRMANN, AND R. ZECKHAUSER (2008): “Betrayal Aversion: Evidence from Brazil, China, Oman, Switzerland, Turkey, and the United States,” *American Economic Review*, 98, 294–310.
- BOHNET, I., B. HERRMANN, AND R. ZECKHAUSER (2010): “Trust and the Reference Points for Trustworthiness in Gulf and Western Countries,” *Quarterly Journal of Economics*, 125, 811–828.
- BOHNET, I. AND R. ZECKHAUSER (2004): “Trust, risk and betrayal,” *Journal of Economic Behavior & Organization*, 55, 467–484.
- BRATMAN, M. (1984): “Two Faces of Intention,” *The Philosophical Review*, 93, 375–405.
- CAMERER, C. AND K. WEIGELT (1988): “Experimental Tests of a Sequential Equilibrium Reputation Model,” *Econometrica*, 56, 1–36.
- CHARNESS, G. (2000): “Responsibility and effort in an experimental labor market,” *Journal of Economic Behavior & Organization, Elsevier*, 42, 375–384.
- (2004): “Attribution and reciprocity in an experimental labor market,” *Journal of Labor Economics*, 22, 665–688.
- CHARNESS, G. AND M. DUFWENBERG (2006): “Promises and partnership,” *Econometrica*, 74, 1579–1601.

- CHARNESS, G. AND D. I. LEVINE (2007): “Intention and stochastic outcomes: An experimental study,” *The Economic Journal*, 117, 1051–1072.
- CHARNESS, G., A. RUSTICHINI, AND J. VAN DE VEN (2013): “Self-confidence and strategic behavior,” Tech. rep., CESifo Working Paper.
- CIKARA, M. AND S. T. FISKE (2013): “Their pain, our pleasure: stereotype content and schadenfreude,” *Annals of the New York Academy of Sciences*, 1299, 52–59.
- CUDDY, A. J. C., S. T. FISKE, V. S. Y. KWAN, P. GLICK, S. DEMOULIN, J.-P. LEYENS, M. H. BOND, J.-C. CROIZET, N. ELLEMERS, E. SLEEBOS, T. T. HTUN, H.-J. KIM, G. MAIO, J. PERRY, K. PETKOVA, V. TODOROV, R. RODRIGUEZ-BAILON, E. MORALES, M. MOYA, M. PALACIOS, V. SMITH, R. PEREZ, J. VALA, AND R. ZIEGLER (2009): “Stereotype content model across cultures: Towards universal similarities and some differences,” *British Journal of Social Psychology*, 48, 1–33(33).
- DOHMEN, T., A. FALK, D. HUFFMAN, AND U. SUNDE (2011): “Individual risk attitudes: Measurement, determinants, and behavioral consequences,” 9, 522–550.
- DUFWENBERG, M. AND G. KIRCHSTEIGER (2004): “A theory of sequential reciprocity,” *Games and Economic Behavior*, 47, 268–298.
- FALK, A., E. FEHR, AND U. FISCHBACHER (2008): “Testing theories of fairness—Intentions matter,” *Games and Economic Behavior*, 62, 287–303.
- FEHR, E., G. KIRCHSTEIGER, AND A. RIEDL (1993): “Does Fairness Prevent Market Clearing,” *Quarterly Journal of Economics*, 108, 437–459.
- FETCHENHAUER, D. AND D. DUNNING (2009): “Do people trust too much or too little?” *Journal of Economic Psychology*, 30, 263–276.
- (2012): “Betrayal aversion versus principled trustfulness: How to explain risk avoidance and risky choices in trust games,” *Journal of Economic Behavior & Organization*, 81, 534–541.
- FISKE, S. T., A. J. CUDDY, AND P. GLICK (2007): “Universal dimensions of social cognition: warmth and competence,” *Trends in Cognitive Sciences*, 11, 77–83.
- FISKE, S. T., A. J. CUDDY, P. GLICK, AND J. XU (2002): “A Model of (Often Mixed) Stereotype Content: Competence and Warmth Respectively Follow From Perceived Status and Competition,” *Journal of Personality and Social Psychology*, 82, 878–902.
- GURDAL, M. Y., J. B. MILLER, AND A. RUSTICHINI (2013): “Why blame?” *Journal of Political Economy*, 121, 1205–1247.
- HUMPHREY, S. J. AND S. MONDORF (2014): “Towards an Understanding of Betrayal Aversion,” Presented at FUR XVI in Rotterdam, 2014.
- KAGEL, J. H. AND K. W. WOLFE (2001): “Tests of fairness models based on equity considerations in a three-person ultimatum game,” *Experimental Economics*, 4, 203–219.
- KASPERSON, R. E., O. RENN, P. SLOVIC, H. S. BROWN, J. EMEL, R. GOBLE, J. X. KASPERSON, AND S. RATICK (1988): “The Social Amplification of Risk: A Conceptual Framework,” *Risk Analysis*, 8, 177–187.

- KNOBE, J. (2006): “The Concept of Intentional Action: A Case Study in the Uses of Folk Psychology,” *Philosophical Studies*, 130, 203–231–.
- LOEWENSTEIN, G. F., E. U. WEBER, C. K. HSEE, AND N. WELCH (2001): “Risk as feelings,” *Psychological Bulletin*, 127, 267–286.
- MELE, A. R. (1992): “Recent Work on Intentional Action,” *American Philosophical Quarterly*, 29, 199–217.
- NERI, C. AND H. ROMMESWINKEL (2014): “Freedom, Power and Interference: An Experiment on Decision Rights,” Available at SSRN: <http://ssrn.com/abstract=2485107> or <http://dx.doi.org/10.2139/ssrn.2485107>.
- NOOTEBOOM, B. (2002): *Trust: Forms, foundations, functions, failures and figures*, Cheltenham, U.K., & Northampton, MA: Edward Elgar Publishing.
- OWENS, D., Z. GROSSMAN, AND R. FACKLER (forthcoming): “The Control Premium: A Preference for Payoff Autonomy,” *AEJ Microeconomics*.
- RABIN, M. (1993): “Incorporating fairness into game theory and economics,” *American Economic Review*, 83, 1281–1302.
- ROSENBERG, S., C. NELSON, AND P. S. VIVEKANANTHAN (1968): “A multidimensional approach to the structure of personality impressions,” *Journal of Personality and Social Psychology*, 9, 283–294.
- SETIYA, K. (2014): “Intention,” in *The Stanford Encyclopedia of Philosophy*, ed. by E. N. Zalta, spring 2014 ed.
- SLOVIC, P. (1987): “Perception of risk,” *Science*, 236, 280–285.
- SLOVIC, P., M. L. FINUCANE, E. PETERS, AND D. G. MACGREGOR (2004): “Risk as Analysis and Risk as Feelings: Some Thoughts about Affect, Reason, Risk and Rationality,” *Risk Analysis*, 24, 311–322.
- WEBER, E. U., A. BLAIS, AND N. E. BETZ (2002): “A domain specific risk attitude scale: Measuring risk perceptions and risk behaviors,” *Journal of Behavioral Decision Making*, 15, 263–290.



## A Appendix: Supplementary Tables & Figures

In Table 5 we report the output of the permutation test for the difference in means, for each between-treatment comparison. The permutation follows the permuted block design of the experiment, where permutations are stratified at the experimental session level. The second column lists the estimated difference between treatments. The third column counts the number of permutations (out of 100,000) where the difference was at least as large as the estimated difference. The fourth column lists the approximate p-values, the proportion of the permuted data where the difference is at least as large as the estimated difference.<sup>27</sup> The 95 percent confidence interval pertains to the p-value, it is a binomial (Clopper-Pearson) confidence interval based on the 100,000 realizations from the permutation scheme.

**Table 5:** (tab: permuteCIs) For each pair-wise between-treatment difference in means, the results of the permutation test are reported below.

Comparison	Difference	Count	P-Value	St. Err.	[95% Conf.	Interval]
AFD vs. xxD	0.08**	4833	0.048	0.001	0.047	0.050
AFD vs. AxD	0.17***	22	0.000	0.000	0.000	0.000
AFD vs. Axx	0.04	18831	0.188	0.001	0.186	0.191
xxD vs. Axx	-0.04	78808	0.788	0.001	0.786	0.791
xxD vs. AxD	0.09**	3253	0.033	0.001	0.031	0.034
Axx vs. AxD	0.13***	381	0.004	0.000	0.003	0.004

100,000 Permutations (player strata)

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$  (one-sided, right)

In Table 6, for the purposes of comparison, we report the outcome of the T-test, where the distribution of mean MAPs in each session is assumed to be normal.

**Table 6:** Below we can see that under the assumption of normality, the p-value of the T-test yields a close approximation of the exact p-value

Comparison	Permutation	T-test
AFD vs. xxD	.0483	.0556
AFD vs. AxD	.0002	.0005
AFD vs. Axx	.1883	.2030
xxD vs. Axx	.7881	.7665
xxD vs. AxD	.0325	.0317
Axx vs. AxD	.0038	.0063

<sup>27</sup>There exists exact p-values for this test, and the Monte-Carlo permutations can approximate them to arbitrary precision.

## B Appendix: Experimental Procedures & Instructions

### B.1 Procedures

**Phase 1** One week before the main experiment sessions (phase 2) 10 students in the graduate program at Bocconi University were invited to take part in an experiment in room “Y”. In a single session, these students participated as the second mover (Person Y) making a pre-commitment of their choice in response to the choice of the first mover (Person X) in every experimental treatment in the following order (1) Treatment Axx, (2) Treatment AxD, (3) Treatment xxD, and finally (4) Treatment AFD. For each treatment they were told that they are Person Y but they are to read the instructions for Person X as their instructions describe every aspect of the game.<sup>28</sup> These students returned after phase 2 of the study to receive their payment.

**Phase 2** Participants were recruited from the Bocconi University online recruitment website administered by Sona-Systems (<http://www.sona-systems.com/>). Each session was given a unique title and description to minimize communication between participants.

When participants arrived they waited until all registered students were present and then were invited into the lab all at once. As they walked in, they selected “code” numbers out of a box and were told immediately: “You have been paired with a another participant”.<sup>29</sup> Subjects were next instructed to seat themselves in the carrel corresponding to their code number. When they were seated the experimenter began with the “Experimenter Script’ presented in Section B.2.

When the script was finished each instruction/decision sheet was folded in half and handed out. In each session there were four sets of instruction/decision sheets, one corresponding to each experimental treatment participants were assigned to. The selection of “code” numbers implemented a permuted block randomization, with a block size of 27 participants and a near-uniform allocation ratio (6,7,7,7).

Participants read the instructions privately and raised their hands to ask clarification questions. When instructions were complete participants filled out a quiz checking their understanding. Next quizzes were collected. Incorrect quizzes were identified by experimental assistants and replaced with a new blank quiz and participants were given an opportunity to ask questions again (the process continued until each participant could demonstrate understanding of the instructions).

After the instruction/decision sheet was collected from each participant they were handed a survey and a receipt form to fill out while the experimenter matched them with Room Y decisions and determined their payment.<sup>30</sup> Next participants were called up one-by-one to be paid based on the choice of the person from Room Y (phase 1) whom they were matched with.

**Phase 3** Students from phase 1 (Room Y) returned one week after phase 2 and were paid for each treatment they participated in. For each treatment their earnings from each participant they were matched with were pooled together. They were paid based on a random selection from their pooled earnings from each treatment. Participants in Room X were not informed that the matching and the payment for participants in Room Y would be conducted precisely in this manner.

---

<sup>28</sup>In Treatment Axx there were no instructions, subjects were simply asked to choose a box.

<sup>29</sup>This was for the random pairing with Room Y participants.

<sup>30</sup>The specific implementation of the matching was not described to the participants of Room X, only that they were uniquely matched with a student from Room Y. The matching was many-to-one and only participants in Room Y were aware of this.

## B.2 Experimenter Script (English Translation): Room X Sessions

1. (Once everyone is seated) Welcome to the study and thank you for participating.
2. First, we ask you to please turn off your mobile devices, not communicate between each other, and leave your desk clear of everything except your student ID and a pen. We will not be using the computers.
3. We will give you a brief overview of the study. It is important that you listen closely. You may ask questions once we have finished reading the instructions (which we will hand to you shortly)
4. This is Room X. When you selected a code number as you walked in this room you were randomly matched with one of the student participants in Room Y. Your identity will be anonymous to them, and theirs to you.<sup>31</sup>
5. In this study you will make a single decision that may influence both your payoffs and the payoffs of the person you are matched with. Please note that there is not a correct or incorrect response, your decision is personal and yours to make.
6. The study will go as follows:
  - (a) You will read the instructions which we will hand to you in a moment.
  - (b) You will answer a short quiz. This quiz will be handed to you just after the instructions are finished. The purpose of the quiz is to confirm that you have perfectly understood the instructions. Be careful, it is important that you answer the key question that you will find on the first page of the instructions *after* you have successfully completed the quiz.
  - (c) Once you have completed the quiz, and after we have checked the correctness of the answers, you can answer the Key Question (the one and only real decision you will make during the course of the experiment!). Keep in mind that there is no relationship between the answers in the quiz and the answer you'll have to give the key question!
  - (d) When everyone has finished we will collect your choices, leave and match your choices with the responses from Room Y, and return.
  - (e) While you wait for us to calculate the amount of your winnings, we will hand out a form asking for your feedback and comments. We will also distribute a survey which is completely anonymous, but there are some personal questions so if you prefer not to respond to some of these, feel absolutely free not to.
  - (f) When we are ready to proceed with payments, after collecting all of the surveys, we will call you one by one to the front of the room. You will need to bring with you the little number that we gave away at the entrance, the completed receipt form (which will be delivered towards the end of the experiment), your student ID, and all your belongings so that you can leave immediately without disturbing others.

---

<sup>31</sup>The Room X and Y designation were chosen so to make it apparent that identities would be kept anonymous. If a participant asked for more details about Room Y students, we responded to specific question, this happened twice. In reality the students in Room Y decided a few days before students in Room X decided and were paid a few days after. We did not reveal this information to keep the saliency of betrayal high.

7. We are nearly ready to have you begin reading the instructions. I would like to emphasize one more time that you read the instructions carefully. This is in your best interest because your earnings from this experiment depend largely on your our answer to the Key Question, the only decision you will make today that has monetary consequences. We remind you, please do not respond to the key question until you have completely read the instructions and responded to the brief quiz.
8. We are now ready, we will give you the instructions, and after a couple of minutes we will hand you the quizzes. Please write your code number (“little number”) at the top of each sheet, so as to avoid confusion with the payments. We will now pause to answer any general questions, please raise your hand and we will come around to you individually.
9. Thank you for your attention, you may begin the instructions. If you have a question at any point, please raise your hand and we will come around to respond.

### B.3 Instructions *Treatment AFD* (English Translation):

**Welcome to the research project! Your code number is: .....**

You are participating in a study in which you will earn some money. The amount you earn will depend on the outcome of a game you will play. At the end of the study, your earnings will be added to your participation fee of 5 Euros, and you will be paid in cash.

*How the study is conducted.* The study is conducted anonymously. Participants will be identified only by code numbers. There is no communication among them. You have been randomly paired with another participant in Room Y, call him/her “**Person Y**”. Person Y will never know your identity and you will never know Person Y’s identity. Your choices are identified solely by your code number and will never be disclosed to anyone. Both you and Person Y are equally informed of these instructions.

*What the study is about.* The study seeks to understand how people decide. You will decide between two alternatives, A and B. Alternative A gives you a certain payoff that does not depend on the choice of Y. Alternative B gives you an outcome that depends on Y’s behavior. Y chooses between options **J** and **K**.

**Payoff Table** The payoff table reads as follows:

Result of your decision	Nature of choice	Your earnings	Earnings of Person Y
Alternative A	Certainty	10	10
Alternative B	Person Y chooses Option J	15	15
	Option K	8	22

The payoff table is as follows

- If you choose A: you and Person Y will each earn 10 Euros.
- If you choose B:
  - If Person Y chooses J, you and Person Y will each earn 15 Euros.
  - If Person Y chooses K, you will earn 8 Euros and Person Y will earn 22 Euros

**KEY QUESTION:** For you to choose Alternative B instead of Alternative A, how large would the probability  $p$  of being paired with a Person Y who chooses option J minimally have to be? (like any probability, it must lie between 0 and 1)

**YOUR ANSWER:** I will choose alternative B for any  $p$  that is at least  ←  
(this means that I choose alternative A for any  $p$  that is less than this cutoff)

*Note: You do not know what the actual value of  $p$  is. Your choice does not influence the value of  $p$ . It is determined by the fraction of persons Y choosing option J. With YOUR ANSWER you indicate how large the fraction of persons Y who choose J has to be before you pick B over A. This is explained in detail on the next page*

### **Conduct of the study C.1.**

1. Before knowing your choice, Person Y has to answer the following question: “Which option, J or K, do you choose in case B?” Everyone will decide in this way. After everyone has decided, we will collect the answer forms. Please fold them so that nobody can see YOUR ANSWER.
2. We will then calculate the percentage of persons Y who chose option J and inform everyone of it. This gives  $p^*$ , the probability of being paired with a Person Y who chose option J.
3. **If  $p^*$  is greater than or equal to your required value of  $p$  (from YOUR ANSWER above), we will follow your instructions. Your earnings will be determined by the choice of the Person Y you are matched with.**
  - (a) If your Person Y chose J, you and Person Y will earn 15 Euros each.
  - (b) If your Person Y chose K, you will earn 8 Euros and Person Y will earn 22 Euros.
4. **If  $p^*$  is less than your required value of  $p$  (from YOUR ANSWER above), we will follow your instructions: You and your Person Y will get the outcome of the certain choice A, namely 10 Euros each.**

### **Completion of Study and Earnings.**

- Before we conduct the study, we ask you to complete a pre-study questionnaire. We will start the study once everyone has correctly filled out this questionnaire.
- You can collect your earnings by presenting your CODE NUMBER FORM at the end of the study. Your earnings will be in an envelope marked with your code number.

## B.4 Instructions *Treatment xxD* (English Translation):

**Welcome to the research project! Your code number is: .....**

You are participating in a study in which you will earn some money. The amount you earn will depend on the outcome of a game you will play. At the end of the study, your earnings will be added to your participation fee of 5 Euros, and you will be paid in cash.

*How the study is conducted.* The study is conducted anonymously. Participants will be identified only by code numbers. There is no communication among them. You have been randomly paired with another participant in Room Y, call him/her “**Person Y**”. Person Y will never know your identity and you will never know Person Y’s identity. Your choices are identified solely by your code number and will never be disclosed to anyone. Both you and Person Y are equally informed of these instructions.

*What the study is about.* The study seeks to understand how people decide. You are confronted with two alternatives, A and B. Alternative A gives you and Person Y a payoff of 10 Euros for sure. Alternative B gives you an outcome that depends on which option (**J** or **K**) is chosen for Y. Each of the 17 cells below contains one symbol, either J or K. The symbols are visible to Person Y but not to you.

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----

Using the online random number service [www.random.org](http://www.random.org), a number between 1 and 17 will be randomly selected for Person Y. If the corresponding cell of the number selected contains J that means option J is selected for Person Y, if the cell contains K that means option K is selected for Person Y.

**Payoff Table** The payoff table reads as follows:

Result of your decision	Nature of choice	Your earnings	Earnings of Person Y
Alternative A	Certainty	10	10
Alternative B	Selection for Person Y Option J	15	15
	Option K	8	22

The payoff table is as follows

- If you choose A: you and Person Y will each earn 10 Euros.
- If you choose B:
  - If option J is selected for Person Y, you and Person Y will each earn 15 Euros.
  - If option K is selected for Person Y, you will earn 8 Euros and Person Y will earn 22 Euros

**KEY QUESTION:** For you to choose Alternative B instead of Alternative A, how large would the probability  $p$  of being paired with a Person Y where option J is selected for them minimally have to be? (like any probability, it must lie between 0 and 1)

**YOUR ANSWER:** I will choose alternative B for any  $p$  that is at least  ←  
(this means that I choose alternative A for any  $p$  that is less than this cutoff)

*Note: You do not know what the actual value of  $p$  is. Your choice does not influence the value of  $p$ . It is determined by the fraction of persons Y who have option J selected for them. With YOUR ANSWER you indicate how large the fraction of persons Y who have option J selected for them has to be before you pick B over A. This is explained in detail on the next page*

### Conduct of the study C.1.

1. After all the options have been selected for those in Room Y, we will first calculate the percentage of people in Room Y who have had option J selected for them and inform everyone of it. This gives  $p^*$ , the probability of being paired with a Person Y who has had option J selected for them.
2. **If  $p^*$  is greater than or equal to your required value of  $p$  (from YOUR ANSWER above), we will follow your instructions. Your earnings will be determined by the option selected for the Person Y you are matched with.**
  - (a) If your Person Y had option J selected for them, you and your Person Y will earn 15 Euros each.
  - (b) If your Person Y had option K selected for them, you will earn 8 Euros and your Person Y will earn 22 Euros.
3. **If  $p^*$  is less than your required value of  $p$  (from YOUR ANSWER above), we will follow your instructions: You and your Person Y will get the outcome of the certain choice A, namely 10 Euros each.**

### Completion of Study and Earnings.

- Before we conduct the study, we ask you to complete a pre-study questionnaire. We will start the study once everyone has correctly filled out this questionnaire.
- You can collect your earnings by presenting your CODE NUMBER FORM at the end of the study. Your earnings will be in an envelope marked with your code number.



## B.5 Instructions *Treatment Axx* (English Translation):

**Welcome to the research project! Your code number is: .....**

You are participating in a study in which you will earn some money. The amount you earn will depend on the outcome of a game you will play. At the end of the study, your earnings will be added to your participation fee of 5 Euros, and you will be paid in cash.

*How the study is conducted.* The study is conducted anonymously. Participants will be identified only by code numbers. There is no communication among them. You have been randomly paired with another participant in Room Y, call him/her “**Person Y**”. Person Y will never know your identity and you will never know Person Y’s identity. Your choices are identified solely by your code number and will never be disclosed to anyone.

*What the study is about.* The study seeks to understand how people decide. You are confronted with two alternatives, A and B. Alternative A gives you a certain payoff that does not depend on the choice of Y. Alternative B gives you an outcome that depends on Y’s behavior. *Persons Y are not aware they are matched with anyone or that any payoffs depend on their behavior.* Each of the 17 cells below contains one symbol, either J or K. The symbols are not visible to you or Person Y.

Without knowing the purpose, Person Y will blindly choose one of the 17 cells below. If the chosen cell contains J, that means Person Y has selected option J. If the chosen cell contains K that means Person Y has selected option K.

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----

Using the online random number service [www.random.org](http://www.random.org), a number between 1 and 17 will be randomly selected for Person Y. If the corresponding cell of the number selected contains J that means option J is selected for Person Y, if the cell contains K that means option K is selected for Person Y.

**Payoff Table** The payoff table reads as follows:

Result of your decision	Nature of choice	Your earnings	Earnings of Person Y
Alternative A	Certainty	10	10
Alternative B	Person Y chooses Option J	15	15
	Option K	8	22

The payoff table is as follows

- If you choose A: you and Person Y will each earn 10 Euros.
- If you choose B:
  - If Person Y chooses the number that corresponds to option J, you and Person Y will each earn 15 Euros.
  - If Person Y chooses the number that corresponds to option K, you will earn 8 Euros and Person Y will earn 22 Euros

**KEY QUESTION:** For you to choose Alternative B instead of Alternative A, how large would the probability  $p$  of being paired with a Person Y who chooses a number that corresponds to option J minimally have to be? (like any probability, it must lie between 0 and 1)

**YOUR ANSWER:** I will choose alternative B for any  $p$  that is at least  ←  
 (this means that I choose alternative A for any  $p$  that is less than this cutoff)

*Note: You do not know what the actual value of  $p$  is. Your choice does not influence the value of  $p$ . It is determined by the fraction of persons Y who chose a number corresponding to option J. With YOUR ANSWER you indicate how large the fraction of persons Y who chose a number corresponding to option J has to be before you pick B over A. This is explained in detail on the next page*

### **Conduct of the study C.1.**

1. After all everyone has made their decision, we will first calculate the percentage of people in Room Y who have chosen a number corresponding to option J and inform everyone in Room X of it. This gives  $p^*$ , the probability of being paired with a Person Y who has chosen a number corresponding to option J.
2. **If  $p^*$  is greater than or equal to your required value of  $p$  (from YOUR ANSWER above), we will follow your instructions. Your earnings will be determined by the the choice of the Person Y you are matched with.**
  - (a) If your Person Y has chosen a number corresponding to option J, you and your Person Y will earn 15 Euros each.
  - (b) If your Person Y has chosen a number corresponding to option K, you will earn 8 Euros and your Person Y will earn 22 Euros.
3. **If  $p^*$  is less than your required value of  $p$  (from YOUR ANSWER above), we will follow your instructions: You and your Person Y will get the outcome of the certain choice A, namely 10 Euros each.**

### **Completion of Study and Earnings.**

- Before we conduct the study, we ask you to complete a pre-study questionnaire. We will start the study once everyone has correctly filled out this questionnaire.
- You can collect your earnings by presenting your CODE NUMBER FORM at the end of the study. Your earnings will be in an envelope marked with your code number.

## B.6 Instructions *Treatment AxD* (English Translation):

**Welcome to the research project! Your code number is: .....**

You are participating in a study in which you will earn some money. The amount you earn will depend on the outcome of a game you will play. At the end of the study, your earnings will be added to your participation fee of 5 Euros, and you will be paid in cash.

*How the study is conducted.* The study is conducted anonymously. Participants will be identified only by code numbers. There is no communication among them. You have been randomly paired with another participant in Room Y, call him/her “**Person Y**”. Person Y will never know your identity and you will never know Person Y’s identity. Your choices are identified solely by your code number and will never be disclosed to anyone. Both you and Person Y are equally informed of these instructions.

*What the study is about.* The study seeks to understand how people decide. You are confronted with two alternatives, A and B. Alternative A gives you a certain payoff that does not depend on the choice of Y. Alternative B gives you an outcome that depends on Y’s behavior. Each of the 17 cells below contains one symbol, either **J** or **K**. The symbols are not visible to you or Person Y.

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----

Person Y will blindly choose one of the 17 cells. If the chosen cell contains J, that means Person Y has selected option J. If the chosen cell contains K that means Person Y has selected option K. This means that regardless of which option Person Y prefers, the option will be selected only if Person Y’s number choice yields that option.

**Payoff Table** The payoff table reads as follows:

Result of your decision	Nature of choice	Your earnings	Earnings of Person Y
Alternative A	Certainty	10	10
Alternative B	Person Y chooses Option J	15	15
	Option K	8	22

The payoff table is as follows

- If you choose A: you and Person Y will each earn 10 Euros.
- If you choose B:
  - If Person Y chooses the number that corresponds to option J, you and Person Y will each earn 15 Euros.
  - If Person Y chooses the number that corresponds to option K, you will earn 8 Euros and Person Y will earn 22 Euros

**KEY QUESTION:** For you to choose Alternative B instead of Alternative A, how large would the probability  $p$  of being paired with a Person Y who chooses a number that corresponds to option J minimally have to be? (like any probability, it must lie between 0 and 1)

**YOUR ANSWER:** I will choose alternative B for any  $p$  that is at least  ←  
(this means that I choose alternative A for any  $p$  that is less than this cutoff)

*Note: You do not know what the actual value of  $p$  is. Your choice does not influence the value of  $p$ . It is determined by the fraction of persons Y who chose a number corresponding to option J. With YOUR ANSWER you indicate how large the fraction of persons Y who chose a number corresponding to option J has to be before you pick B over A. This is explained in detail on the next page*

### **Conduct of the study C.1.**

1. After all everyone has made their decision, we will first calculate the percentage of people in Room Y who have chosen a number corresponding to option J and inform everyone in Room X of it. This gives  $p^*$ , the probability of being paired with a Person Y who has chosen a number corresponding to option J.
2. **If  $p^*$  is greater than or equal to your required value of  $p$  (from YOUR ANSWER above), we will follow your instructions. Your earnings will be determined by the the choice of the Person Y you are matched with.**
  - (a) If your Person Y has chosen a number corresponding to option J, you and your Person Y will earn 15 Euros each.
  - (b) If your Person Y has chosen a number corresponding to option K, you will earn 8 Euros and your Person Y will earn 22 Euros.
3. **If  $p^*$  is less than your required value of  $p$  (from YOUR ANSWER above), we will follow your instructions: You and your Person Y will get the outcome of the certain choice A, namely 10 Euros each.**

### **Completion of Study and Earnings.**

- Before we conduct the study, we ask you to complete a pre-study questionnaire. We will start the study once everyone has correctly filled out this questionnaire.
- You can collect your earnings by presenting your CODE NUMBER FORM at the end of the study. Your earnings will be in an envelope marked with your code number.