

Discussion Paper Series – CRC TR 224

Discussion Paper No. 223  
Project A 01, B 03

Fishing for Good News:  
Motivated Information Acquisition

Si Chen <sup>1</sup>  
Carl Heese <sup>2</sup>

August 2021  
(First version : October 2020)

<sup>1</sup> Corresponding author; Address: Department of Economics, University of Vienna, Oskar-Morgenstern-Platz 1, 1090 Vienna, Austria; Email: [si.chen@univie.ac.at](mailto:si.chen@univie.ac.at)

<sup>2</sup> Address: Department of Economics, University of Vienna, Oskar-Morgenstern-Platz 1, 1090 Vienna, Austria; Email: [carl.heese@univie.ac.at](mailto:carl.heese@univie.ac.at)

Funding by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) through CRC TR 224 is gratefully acknowledged.

# Fishing for Good News: Motivated Information Acquisition\*

Si Chen<sup>†</sup>

Carl Heese<sup>‡</sup>

July 28, 2021

## Abstract

The literature on motivated reasoning argues that people skew their beliefs to feel moral when acting selfishly. We study information acquisition of decision-makers with a motive to form positive moral self-views and a motive to act selfishly. Theoretically and experimentally, we find that a motive to act selfishly makes individuals ‘fish for good news’: they are more likely to continue (stop) acquiring information, having received mostly information suggesting that acting selfishly is harmful (harmless) to others. We find that fishing for good news may improve social welfare. Finally, more intelligent individuals have a higher tendency to fish for good news.

In this paper, we study empirically and theoretically the information acquisition of a decision-maker for whom information might reconcile two motives that govern her utility. In many situations, decisions are guided both by an

---

\*For helpful discussions and comments, the authors are grateful to Thomas Dohmen, Lorenz Götte, Lena Janys, Botond Köszegi, Sebastian Kube, George Loewenstein, Florian Zimmermann, as well as audience at Bonn Applied Micro Workshop, Workshop on Belief-Dependent Preferences (Copenhagen), Workshop on Behavioral Game Theory (Norwich), Workshop on Contracts and Incentives (Munich) and 10th them on Social Preferences (Konstanz). Funding by the German Research Foundation (DFG) through CRC TR 224, Project A01 (Chen) and Project B03 (Heese) is gratefully acknowledged.

<sup>†</sup>Corresponding author; Address: Department of Economics, University of Vienna, Oskar-Morgenstern-Platz 1, 1090 Vienna, Austria; Email: si.chen@univie.ac.at

<sup>‡</sup>Address: Department of Economics, University of Vienna, Oskar-Morgenstern-Platz 1, 1090 Vienna, Austria; Email: carl.heese@univie.ac.at

egoistic motive—a desire to maximize personal gains—and a moral motive. Growing empirical evidence shows that the moral motive is often belief-based: people want to ‘feel moral’, whether their decisions are actually moral or not (for reviews, see Kunda, 1990; Bénabou and Tirole, 2006; Gino *et al.*, 2016). This motive to feel moral might compete with the individual’s egoistic motive if she believes that maximizing her personal gain is detrimental for others. That is, she cannot behave selfishly while feeling moral.

However, individuals are sometimes uncertain about whether a self-benefiting choice is harmful to others. Under uncertainty, new information brings the chance of reconciling the egoistic and the moral motive since it may suggest that an egoistic decision is also moral.

Examples are numerous where an egoistic motive might conflict with the motive to feel moral: doctors receive commissions for prescribing certain drugs, but prescribing the commissioned drug might harm the patients’ health. A human resource manager may have personal preferences for job candidates of certain ethnicity or gender, but hiring decisions based on her personal taste might harm the company’s performance and the job candidates’ careers. Consumers might find it economical to purchase fast fashion products, but doing so might support unethical production.

In these situations, systematic biases in information collection may affect the social outcome of decisions. How doctors gather information about the patients’ medical needs may affect the suitability of the prescribed drugs. When a human resource manager’s personal taste against minority job candidates sways how she informs herself about the candidates’ job-related qualities, the hiring outcome is biased. What information consumers acquire about production conditions of goods may affect the prevalence of unethical production.

In this paper, we theoretically and experimentally study how the potential trade-off between an egoistic and a moral motive shapes the individuals’ information acquisition and the welfare consequences. To do so, we consider environments with ‘rich’ information sources. This ensures that the predicted and the observed information choices are not driven by arbitrary exogenous limitations imposed by the study. It allows us to observe uncensored data on the individuals’ information strategies and uncover novel phenomena. Further, we consider a baseline in which the egoistic motive is removed from the

decision. Comparing this baseline with the scenario with an egoistic motive, we can study the causal effects of the trade-off between competing motives.

We find that having an egoistic motive in the decision systematically biases the information acquisition. Theoretically, we show that it causes individuals to ‘*fish for good news*’: the central feature of optimal strategies is that individuals are more likely to *continue* acquiring information after having received mostly bad news, and more likely to *stop* after having received mostly good news. Here, ‘bad news’ is a piece of information indicating that behaving selfishly harms the other, and ‘good news’ indicate the opposite. In a laboratory experiment, we find that individuals do fish for good news.

Surprisingly, both in theory and in the experiment, we find that *fishing for good news* may *reduce* the harm that the decisions cause on others. The intuition is that having a self-benefiting option causes individuals to fish for good news, which may spur information acquisition overall. In turn, individuals might make better-informed decisions and cause less harm on others.

Further, we find that the tendency to *fish for good news* is stronger among more intelligent individuals—evidence that fishing for good news is more likely a strategic behaviour than a result of cognitive limitations.

To guide the empirical analysis, we formulate a theoretical framework with two features: first, our goal is to analyze how individuals make a trade-off between the desire to behave selfishly and the desire to feel moral. For this, we consider an agent with two motives—she gains utility from her material payoff and from the *belief* about her choice’ being harmless to others.<sup>1</sup> Second, the predicted behaviour should *only* be driven by the trade-off between the two motives and not constrained or confounded by other factors. Therefore, we consider an environment without *any* exogenous restriction on the agent’s choice of information and where information comes at no cost: information about which option harms others arrives continuously, and at every point of time, the agent can decide to either stop or continue observing the incoming information.<sup>2</sup>

---

<sup>1</sup>The study of belief-based utility has a long tradition in economics (e.g., Loewenstein, 1987; Geanakoplos *et al.*, 1989; Akerlof and Kranton, 2000; Köszegi, 2006).

<sup>2</sup>We study a dynamic setup because it allows us to implement an intuitive and ‘rich’ information environment in the laboratory. Consistent with this model setup, we ask the subjects in the experiment to decide whether to stop or continue receiving information, without time limit.

The model brings forward the following intuition why individuals fish for good news: after bad news indicating that a materially self-benefiting option is likely to harm others, an individual may be inclined not to choose this option to avoid a low belief-based utility. Then, more information comes in handy. First, if good news arrived, it may revert her decision, so that she chooses the self-benefiting option instead. Second, even if bad news arrived and she decided to forgo the self-benefit, she would be more certain that doing so actually spares the other from harm. Either way, she is better-off acquiring further information. In contrast, when the individual has received mostly good news indicating that the self-benefiting option is likely harmless, she may be inclined to capture the self-benefits. Then, collecting further information bears the risk that the self-benefiting option becomes morally unacceptable, an effect discouraging her from acquiring further information.

To empirically investigate the information acquisition, we conduct a laboratory experiment. The controlled laboratory environment allows us to address three challenges facing an empirical investigation of information acquisition. First, exogenous variation in the motives is required to pin down the causal effect of having two potentially conflicting motives on information acquisition. Second, when the available information strategies are limited, this leads to censored data, confounding the analysis. Third, individuals' heterogeneous prior beliefs, access to information and interpretation of the information can confound the observed information acquisition strategy.

Our experiment addresses these challenges. In the experiment, we induce the moral motive by having the subjects make a binary dictator decision. In the dictator decision, one of the two options reduces the payoff of a receiver, while the other does not. The dictator does not know which option is harmful to the receiver. We fix the dictator's prior belief about the likelihood of each option being the harmful one. We exogenously vary the existence of an egoistic motive by randomly assigning the dictators into two treatments. In one treatment, one option increases the dictator's own payment, and she knows which option is self-benefiting. Thus, the dictator has an egoistic motive to choose the self-benefiting option. In the other treatment, the dictator's payment is not at stake in the dictator decision. This treatment serves as a baseline. Before making the decision, the dictator can acquire information about which option harms the receiver. Information comes in pieces, is free,

and the dictator can stop or continue receiving information at anytime. The information has a clear Bayesian interpretation, and we provide the dictators with the Bayesian posterior beliefs after each piece of information.

This dynamic setup offers the rich information access that is necessary for analysing the dictators' optimal information acquisition strategies. In contrast, a static setup, where dictators can decide whether to receive one piece of information or not, cannot be used to analyze optimal information strategies. We discuss this point in detail in Section 1.6.4.

This paper contributes to the understanding of the formation of motivated beliefs. Previous studies find that individuals react to *exogenous* information in a self-serving manner (Eil and Rao, 2011; Mobius *et al.*, 2011; Gneezy *et al.*, 2016; Falk and Szech, 2016; Exley and Kessler, 2018), have selective *memories* (Zimmermann, 2020), and directly manipulate their beliefs in settings *without information* (Di Tella *et al.*, 2015; Haisley and Weber, 2010). In the psychology literature, Ditto and Lopez (1992) document that individuals require less supportive information to reach their preferred conclusion, possibly due to individuals' over-interpreting preferred information. In comparison, we facilitate Bayesian updating in the experiment and focus on the individuals' use of information acquisition *per se* as an instrument to form motivated beliefs, rather than the fact that information deemed more valid leads to a conclusion faster.

In the motivated beliefs literature, the previous work on information choices has focused on one-shot information that reveals whether acting selfishly harms others. The existing finding is that individuals avoid information and stay willfully ignorant (Dana *et al.*, 2007; Feiler, 2014; Grossman, 2014; Golman *et al.*, 2017; Serra-Garcia and Szech, 2019). Willful ignorance is an important phenomenon in many applications. In this paper, we examine the effect of the trade-off between an egoistic motive and a moral motive on information acquisition. To achieve this end, the individuals' choices of information should be driven by this trade-off and not by exogenous restrictions on their choices imposed by the researchers. Any exogenous restriction, e.g., the restriction to a binary information choice between being fully informed and not informed, may censor the observed information choices and hence confound the analysis. Therefore, we consider rich information environments. Empirically, this allows

us to unmask the following information acquisition strategy: individuals *fish for good news* instead of avoiding information.

We find that *fishing for good news* increases social welfare in our data, an effect that contrasts the negative connotation of willful ignorance. Notably, the identification of the welfare effect is enabled by our baseline treatment, where we remove the egoistic motive. This design contrasts the existing literature, which typically compares the full information scenario and a ‘hidden information’ scenario.<sup>3</sup>

Fishing for good news is predicted as the optimal strategy in our novel model with rich information choices. The model also connects to the theoretical literature on willful ignorance. The key feature of our model is that some beliefs are more desirable than others. This feature is shared by several psychological factors that this literature has brought forward as explanations for willful ignorance—most prominently self-image concerns (Grossman and van der Weele, 2017) and moral constraints as in Rabin (1995). We discuss these existing models and another interpretation of belief-based utility, guilt aversion (Battigalli and Dufwenberg, 2007), in the context of our model in Section 1.7: the general prediction is that the trade-off between the belief-based motive and the egoistic motive causes individuals to fish for good news. Apart from that, when restricting agents to the binary choice between full and no information, the model predicts the avoidance of information.

Our paper joins an emerging literature leveraging tools from the theory of *interpersonal* strategic information transmission to study motivated beliefs. For example, Bénabou and Tirole (2006) and Bodner and Prelec (2003) use a signaling game to model self-image concerns. Grossman and van der Weele (2017) build upon this framework and study willful ignorance. Connecting to the literature on the disclosure of hard information, Hagenbach and Koessler (2021) propose a model of hard evidence to study selective memory.

Our paper connects to the literature on Bayesian persuasion (Kamenica and Gentzkow, 2011), which studies information transmission from one individual to another when the latter does not share the same interests as the former. We show that there is an analogy between the model in this paper

---

<sup>3</sup>In a later paper, Exley and Kessler (2021) implement a baseline without egoistic motive to study the motive behind information avoidance. They provide no welfare analysis.

and models of *interpersonal* persuasion: a decision-maker who acquires information to align two internal motives faces a similar situation as a person who sends information to another individual to align the other individual's belief (and his subsequent choice) with her own interests (see Section 4.1 for a detailed discussion). This analogy allows us to use techniques from existing models of Bayesian persuasion to study motivated information acquisition.<sup>4</sup>

We organize the rest of the paper as follows: In Section 1, we present the theoretical results and testable predictions. In Section 2, we detail the experimental design. In Section 3, we present the empirical analysis of information acquisition. In Section 4, we provide further discussion on the theoretical and the empirical analyses. In Section 5, we conclude and propose ideas for future research.

## 1 Theory

We propose a formal model to analyze an agent's information acquisition in a decision where she has an egoistic motive and a motive to believe that her decision is moral. To highlight the effect of the egoistic motive, we also study the scenario in which the egoistic motive is removed.

After stating the main result, we make three points that are important to understand the interplay between the egoistic and the moral motive. We then move on to connect these three points and lay out the formal proof of the main result. Afterwards, we discuss interpretations of the belief-based utility; in particular, we relate them to existing models in the literature. Finally, we derive five testable predictions that guide our empirical investigation.

### 1.1 A model of conflicting motives

An agent (she) has to make a decision between two options  $x$  and  $y$ . There is an unknown binary state  $\omega \in \{X, Y\} = \Omega$  and the prior belief is that the probability of  $X$  is  $\Pr(X) \in (0, 1)$ . A passive agent, whom we hereafter refer

---

<sup>4</sup>In *intra-personal* persuasion, the commitment assumption of Bayesian persuasion entails that the agent cannot hide the acquired information from herself, which does not demand a high level of sophistication.

to as *the other* (he), can be affected by the agent's decision between  $x$  and  $y$ . When the agent chooses an option that does not match the state, i.e.,  $x$  in  $Y$  or  $y$  in  $X$ , the option has a negative externality of  $-1$  on the other and otherwise not.

**Preferences.** The agent's preferences are governed by two motives. First, if choosing  $x$ , the agent receives a state-independent *remuneration*  $r \geq 0$ , while she receives no remuneration if choosing  $y$ . When  $r > 0$ , the remuneration constitutes an egoistic motive to choose  $x$ .<sup>5</sup> The case  $r = 0$  serves as the benchmark without egoistic motive. Second, the agent has a moral motive. She dislikes the belief that her decision harms the other. We model this as the agent receiving a utility  $u(a, q)$  when she believes that her choice  $a$  is harmless for the other agent with probability  $q$ . We discuss different interpretations of the belief-based utility in detail in Section 1.7. Note that, when the agent believes that state  $X$  holds with probability  $p$ , she believes that  $x$  is harmless with probability  $q = p$  and that  $y$  is harmless with probability  $q' = 1 - p$ . Then, when she chooses  $a \in \{x, y\}$ , her utility is given by

$$U(a, p; r) = \begin{cases} u(a, p) + r & \text{if } a = x, \\ u(a, 1 - p) & \text{if } a = y. \end{cases} \quad (1)$$

The *belief-based utility*  $u$  is weakly increasing in the second argument. We let  $u(x, 1) = u(y, 1) = 0$ . That is, the dictator feels no disutility if she is certain that her choice does not harm the other. We also call  $u$  the (*preference*) *type* of the agent. For concreteness, below we state a parametric example of the belief-based utility  $u$  to facilitate the interpretation of what a type is.

**A parametric example of the type.** Consider

$$u(a, q) = -\theta(1 - q)^2. \quad (2)$$

Here,  $\theta = u(a, 0)$  is the disutility from choosing an action  $a$  that harms the other with certainty, i.e., when  $q = 0$ .

---

<sup>5</sup>The remuneration here is a token standing for not only monetary interests but also any private interest that the agent might have. In the example of a discriminatory human resource manager, the private interest can be the utility of her choosing a candidate of her personally preferred gender.

**Information and strategies.** Before deciding between  $x$  and  $y$ , the agent can acquire information about the state at no cost. Let  $\mu(\omega) = 0$  if  $\omega = X$  and  $\mu(\omega) = 1$  if  $\omega = Y$ . Time is continuous and at every instant in time the agent can observe an *information process*  $(Z_t)_{t \geq 0}$  given by  $dZ_t = \mu(\omega)dt + dW_t$  where  $(W_t)_{t \geq 0}$  is a standard Brownian motion. The posterior probability that the agent assigns to the state  $X$  at the time  $t$  is

$$p_t = \Pr(\omega = X | (Z_s)_{s \leq t}).$$

At every point of time, the agent can decide to either stop or continue observing the process  $(Z_t)_{t \geq 0}$ , depending on the information she has already received. When the agent stops at  $t \geq 0$ , subsequently the agent chooses an action  $a$  that maximizes her payoffs, i.e.,  $a \in \max_{a \in \{x, y\}} U(a, p_t; r)$  and the game ends. Formally, a strategy of the agent is a real-valued stopping time  $\tau$  adapted to the natural filtration generated by the information process.

For technical reasons, we impose the ‘coarseness condition’ that the agent stops and takes a decision when  $p_t \leq \epsilon$  or  $p_t \geq 1 - \epsilon$ , for some positive, but *arbitrarily* small  $\epsilon \approx 0$ .<sup>6</sup> This is to rule out strategies where the agent observes the information process infinitely with positive probability.

## 1.2 Equilibrium characterization

**Lemma 1** *There are cutoffs  $p_l \leq p_0 \leq p_h$ , so that the following constitutes a subgame perfect equilibrium: the agent continues to observe the information process as long as  $p_l < p_t < p_h$ , and stops whenever  $p_t \leq p_l$  or  $p_t \geq p_h$ .*

The proof is in Appendix B. Lemma 1 shows the existence of an equilibrium. To show the lemma, we leverage an insight from the analysis of Bayesian persuasion (Kamenica and Gentzkow, 2011).

Since there are no cost of observing the information process, any Nash equilibrium must maximize  $E(V(p_\tau))$  with  $V(p) = \max_{a \in \{x, y\}} U(a, p; r)$  and where  $p_\tau$  is the stopped belief. This implies that all equilibria are payoff equivalent.

---

<sup>6</sup>In the experiment, posteriors are rounded to two decimal places, so that e.g beliefs below 1% are identified with certainty, essentially implementing  $\epsilon = 0.01$ .

**Lemma 2** *There is a unique subgame perfect equilibrium in which the agent stops observing the information process whenever he is indifferent between stopping and continuing.*

We prove Lemma 2 in Appendix B and show that the equilibrium in Lemma 2 is given by the belief cutoffs  $p_l$  and  $p_h$  as follows: let  $\bar{V}$  be the smallest concave function with  $\bar{V}(p) \geq V(p)$  for all  $p \in [\epsilon, 1 - \epsilon]$ . If  $\bar{V}(p_0) = V(p_0)$ , then  $p_h = p_l$ . Otherwise,  $I = (p_l, p_h)$  is the largest open interval in  $[\epsilon, 1 - \epsilon]$  with  $\bar{V}(p) > V(p)$  for all  $p \in I$ .

For the ease of exposition, we focus on the equilibrium in Lemma 2 and simply use  $p_l$  and  $p_h$  to refer to it. One can show that the main result does not depend on this equilibrium selection.

### 1.3 Result: fishing for good news

The key difference between the scenarios with and without an egoistic motive is that when there is an egoistic motive ( $r > 0$ ), the agent makes a trade-off between the desire for the remuneration and a desire for accurate beliefs. In this section, we analyze how this trade-off affects the agent's behaviour. Our main result, Theorem 1, concerns the effect on the *intensive margin* of information acquisition, i.e., the agent's decision to continue or stop acquiring information once she has started. In Appendix D.1.1, we discuss the *extensive margin* of information acquisition, i.e., the agent's decision whether to acquire *any* information.

In Theorem 1, we consider all types that plan on acquiring some information and use it in a 'responsive' way, i.e., choosing  $y$  after information indicating that  $y$  is harmless to the other, and  $x$  after information indicating that  $x$  is harmless to the other.<sup>7</sup> The theorem shows that, when  $r > 0$ , the agent stops and chooses  $y$  only at a more extreme belief in  $y$  being harmless, i.e.,  $1 - p_l(r) \geq 1 - p_l(0)$ . Conversely, the agent is willing to stop and choose  $x$  at a less extreme belief in  $x$  being harmless, i.e.,  $p_l(r) \leq p_l(0)$ .

---

<sup>7</sup>The theoretical analysis in this section focuses on the trade-off between belief-based utility and the remuneration. It turns out that, when  $r > 0$ , there are types who choose to acquire some information but choose  $x$  regardless of the information they receive. However, such behaviour is not driven by a meaningful trade-off between belief-based utility and remuneration.

**Theorem 1 (Fishing for good news)** *Take any preference type  $u$  and let  $\bar{r} > 0$ . If it is strictly optimal in equilibrium to choose  $y$  at  $p_l(r)$  and  $x$  at  $p_h(r)$  when  $r = \bar{r}$  and also when  $r = 0$ , then*

$$p_h(\bar{r}) \leq p_h(0), \quad (3)$$

$$1 - p_l(\bar{r}) \geq 1 - p_l(0). \quad (4)$$

Theorem 1 reveals an asymmetry. In intuitive terms, (3) shows that to convince herself to choose the remunerative option  $x$ , the agent needs less information supporting the innocuousness of  $x$  (good news). (4) shows that for choosing the non-remunerative option  $y$  the agent needs more information opposing the innocuousness of  $x$  (bad news). Taken together, the agent ‘fishes for good news’ to choose the remunerative option.

We elaborate the intuition for such behaviour below in Section 1.4 and lay out the proof of Theorem 1 in Section 1.5.

## 1.4 Intuition for fishing for good news

In Section 1.4.1 to 1.4.3, we highlight three points that are important to understand Theorem 1. In Section 1.4.4, we describe the agent’s information acquisition strategy in the scenario without egoistic motive, that is when both options are not remunerative.

In each of the subsections 1.4.1 to 1.4.3, we first discuss the point in intuitive terms, and then state a formal result.

### 1.4.1 The desired belief

The first point is that, when one option is remunerative, the agent prefers higher beliefs in the state where this option is harmless. This is because, when believing that the remunerative option is harmless to the other, she can capture the reward without having a bad conscience. In contrast, when she believes that the remunerative option is harmful to the other, she has to make a trade-off between a clear conscience and the remuneration—there is a *moral dilemma*.

Formally, let  $r > 0$ . Recall that in equilibrium, the agent eventually either

stops at  $p_h$  or  $p_l$ , and that at  $p_h$  she has a higher belief about the likelihood that  $x$  is harmless. Similar to the intuition sketched in the previous paragraph, the following result shows that, in equilibrium, the agent is better off when she stops at the higher belief  $p_h$  compared to when she stops at  $p_l$ . The proof is in the Appendix.

**Lemma 3** *For all  $r > 0$ , the agent chooses  $x$  when stopping at the belief  $p_t = p_h$ . Further, if the agent weakly prefers to choose  $y$  when holding the belief  $p_t = p_l$ , then,  $V(p_l) < V(p_h)$ .*

In the next two subsections, we analyze at which beliefs the agent stops. In Section 1.4.2, we analyze at which belief  $p_l$  the agent stops and chooses  $y$ . In Section 1.4.3, we analyze at which belief  $p_h$  the agent stops and chooses  $x$ .

#### 1.4.2 Waiting for good news

In the previous subsection we made the point that the agent prefers to believe that the likelihood of the remunerative option being harmless is high. The second point is: when she believes this likelihood to be low so that she is inclined to choose the non-remunerative option, the agent prefers to continue observing the arriving information. One intuitive reason for this behaviour is that she hopes to receive ‘good news’ so that her belief increases, making it optimal to choose the remunerative option. The second reason is that even if no good news arrived, her belief in the innocuousness of the non-remunerative option would increase, and so would her belief-based utility when choosing it. In any case, she is better off continuing. Formally, we show the following result.

**Lemma 4** *For all  $r > 0$ : if the agent weakly prefers to choose  $y$  when holding the belief  $p_t = p_l$ , then  $p_l = \epsilon \approx 0$ .*

**Proof.** At each point of time  $t \geq 0$ , the equilibrium strategy  $\tau^*$ , given by  $p_l$  and  $p_h$ , maximizes the continuation payoff  $E(V(p_\tau)|(Z_s)_{s \leq t})$ ,

$$\begin{aligned} E(V(p_{\tau^*})|(Z_s)_{s \leq t}) &= \frac{p_h - p_t}{p_h - p_l} V(p_l) + \frac{p_t - p_l}{p_h - p_l} V(p_h) \\ &= u(y, 1 - p_l) + \frac{p_t - p_l}{p_h - p_l} [V(p_h) - V(p_l)], \end{aligned} \quad (5)$$

where, for the first equality, we used that  $E(p_{\tau^*} | (Z_s)_{s \leq t}) = p_t$  by Bayes-consistency.<sup>8</sup> For the second equality, we used that the agent chooses  $y$  at  $p_l$ , so that  $V(p_l) = u(y, 1 - p_l)$ . We see that the continuation payoff strictly decreases in  $p_l$  since the likelihood of reaching  $p_h$ , that is  $\frac{p_t - p_l}{p_h - p_l}$ , decreases in  $p_l$  and since the utility  $u(y, 1 - p_l)$  when reaching the lower belief  $p_l$ , also decreases in  $p_l$ . We conclude that, unless the agent is certain that  $y$  is harmless, she would like to continue observing the arriving information, thus,  $p_l = \epsilon$ .<sup>9</sup>

■

### 1.4.3 Good enough news

The third point is that when she believes that the remunerative option is likely to be harmless, then she decides if to stop and choose this option by making a trade-off between her belief-based utility with the remuneration: on the one hand, if she continues, her belief in this option being harmless may increase further, allowing her to have a better conscience when choosing it. However, continuing to acquire information bears the risk of observing information that makes the remunerative option unacceptable, i.e., that leads her to update to a low posterior and to choose the non-remunerative option.

Formally, at each point of time  $t \geq 0$ , the equilibrium strategy  $\tau^*$ , given by  $p_l$  and  $p_h$ , maximizes the continuation payoff  $E(V(p_\tau) | (Z_s)_{s \leq t})$ . From Lemma 4, we take  $p_l \approx 0$ , so that  $E(V(p_\tau) | (Z_s)_{s \leq t}) \approx \Pr(p_\tau = p_h | (Z_s)_{s \leq t})V(p_h)$ . For expositional purposes only, let  $u(x, p)$  be continuously differentiable for  $q > p_0$ . Using  $V(p_h) = u(x, p_h) + r$ , the first-order condition with respect to  $p_h$  is

$$0 = \Pr(p_\tau = p_h | (Z_s)_{s \leq t}) \frac{\partial u(x, p_h)}{\partial p_h} + \frac{\partial \Pr(p_\tau = p_h | (Z_s)_{s \leq t})}{\partial p_h} (u(x, p_h) + r), \quad (6)$$

which shows that the agent makes a trade-off between the marginal increase in belief-based utility from stopping at a higher belief  $p_h$  and the marginal decrease in the likelihood of stopping at  $p_h$ , which comes with the remuneration

---

<sup>8</sup>Given the strategy  $\tau^*$ , Bayes-consistency implies  $\Pr(p_{\tau^*} = p_h | (Z_s)_{s \leq t}) = \frac{p_t - p_l}{p_h - p_l}$  and  $\Pr(p_{\tau^*} = p_l | (Z_s)_{s \leq t}) = \frac{p_h - p_t}{p_h - p_l}$ .

<sup>9</sup>Recall the technical restriction that the agent has to stop if  $p_t = \epsilon$  where  $\epsilon \approx 0$  is arbitrarily small.

$r$ . Rewriting (6),<sup>10</sup>

$$0 = \frac{\partial u(x, p_h)}{\partial p_h} p_h - (u(x, p_h) + r). \quad (7)$$

Recalling  $u(x, 1) = 0$ , one sees from (7) that when the marginal increase in belief-based utility is relatively small for high beliefs  $p_t \approx 1$ , precisely when  $\frac{\partial u(x, 1)}{\partial p_h} < r$ , the agent is willing to stop and choose  $x$  before she is certain that  $x$  is harmless. One may say that the agent stops when she has received ‘good enough news’.

#### 1.4.4 The scenario without egoistic motive

Next, we turn to the agent’s decision between stopping and continuing acquiring information in the benchmark scenario where both options are *not* remunerative. The agent’s utility depends solely on her belief about the likelihood that her action does not harm the other. In this scenario, the agent stops acquiring information only when further certainty no longer increases her utility.

For the agent types  $u$  with  $u' > 0$ , the more certain they are that their decision does not harm the other, the higher their utility would be. For these agent types, it is optimal to acquire as much information as possible. Other agent types have a threshold level of certainty. They are content when believing that it is sufficiently likely that they can spare the other from harm. Any further certainty beyond the threshold does not increase their belief-based utility. At the threshold, such types are indifferent between continuing and stopping, so they may as well stop. This behaviour mirrors that of *satisficing* as in Simon (1955).<sup>11</sup>

Formally, the threshold level of certainty is

$$l(a) = \min \{q : u(a, q) = 0\},$$

which implies that  $u(a, q) = 0$  for all  $q \geq l$ , recalling that  $\max_q u(a, q) = 0$  and

---

<sup>10</sup>Recall that  $\Pr(p_{\tau^*} = p_h | (Z_s)_{s \leq t}) = \frac{p_t - p_l}{p_h - p_l} \approx \frac{p_t}{p_h}$ , so that  $\frac{\partial \Pr(p_{\tau^*} = p_h | (Z_s)_{s \leq t})}{\partial p_h} \approx -\frac{p_t}{p_h^2}$ . Plugging this into (6) gives  $\frac{p_t}{p_h} \frac{\partial u(x, p_h)}{\partial p_h} p_h - \frac{p_t}{p_h^2} (u(x, p_h) + r) = 0$ , which simplifies to (7).

<sup>11</sup>We are unaware of a literature where the satisficing behaviour concerns beliefs instead of outcomes.

that  $u$  is increasing.<sup>12</sup> The following result formally describes the equilibrium when  $r = 0$ . It shows that the agent acquires information until she reaches her threshold level of certainty  $l(a)$ , unless  $l(a) > 1 - \epsilon$ .<sup>13</sup> Here, it may be the case that the threshold level of certainty is already reached at the prior for one of the options, so that she stops directly.

**Lemma 5** *Let  $r = 0$ . If  $\max_{a \in \{x, y\}} l(a) \leq 0.5$  or  $p_0 \in (1 - l(y), l(x))^c$ , then  $p_l = p_0 = p_h$ . If  $p_0 \in [1 - l(y), l(x)]$ , then  $p_l = \max\{\epsilon, 1 - l(y)\}$  and  $p_h = \min\{l(x), 1 - \epsilon\}$ .*

## 1.5 Proof of Theorem 1

Take any ‘responsive type’  $u$ , meaning that it is strictly optimal for the type to choose the option  $y$  at the belief  $p_l$  and the option  $x$  at the belief  $p_h$  when  $r > 0$  and when  $r = 0$ . First, recall from Lemma 4 that  $p_l(r) = \epsilon$ . Hence  $p_l(r) \leq p_l(0)$ , which shows the second part of Theorem 1, (4).

To show the first part of Theorem 1, (3), first we note that it follows from Lemma 5 that  $p_h(0) \in \{1 - \epsilon, l(x), p_l(0)\}$ . When  $p_h(0) = p_l(0)$ , the agent is not responsive, so the precondition of the theorem is not fulfilled. It remains to establish that  $p_h(r) \leq \min\{1 - \epsilon, l(x)\}$ . Clearly  $p_h(r) \leq 1 - \epsilon$  since the agent has to stop at  $1 - \epsilon$  necessarily. Finally, we show that  $p_h(r) \leq l(x)$ . Given the definition of  $l(x)$  in (8), we know that either  $l(x) = 1$  or  $\frac{\partial u(x, p)}{\partial p} = 0$  for all  $p > l(x)$ .<sup>14</sup> If  $l(x) = 1$ , clearly  $p_h(r) \leq l(x)$ . For the second case, observe that the derivative of the objective function with respect to  $p_h(r)$ , which is the left hand side of (7), is strictly negative for any  $p > l(x)$  when  $r > 0$ . This follows since  $\frac{\partial u(x, p)}{\partial p} = u(x, p) = 0$  for all  $p > l(x)$ . Hence,  $p_h(r) \leq l(x)$ . This finishes the proof of the claim that  $p_h(r) \leq \min\{1 - \epsilon, l(x)\}$ , and thereby the proof of Theorem 1.

<sup>12</sup>The discussion in Section 1.7 explains how the threshold level may e.g. be interpreted as a *moral constraint* (Rabin, 1995).

<sup>13</sup>Recall that, for technical reasons, we restrict the agent’s strategies, imposing that the agent has to stop at  $p_t = \epsilon$  and  $p_t = 1 - \epsilon$  for  $\epsilon \approx 0$  arbitrarily small.

<sup>14</sup>In particular,  $l(x) < 1$  implies the continuous differentiability of  $u(x, p)$  for  $p > l(x)$ .

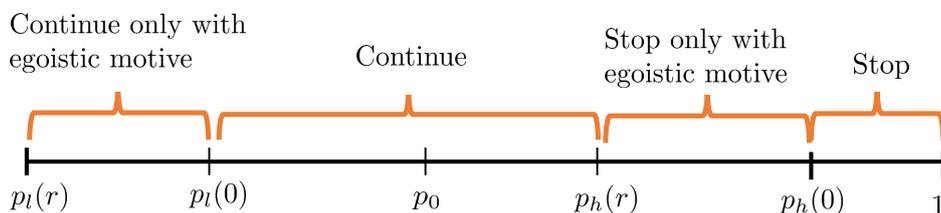
## 1.6 Testable predictions

The following testable predictions will guide our experimental study on the information acquisition presented in Section 2 and 3.

### 1.6.1 People fish for good news

Individuals stop acquiring information when their beliefs reach either the lower cutoff  $p_l$  or the upper cutoff  $p_h$  (see Lemma 1). Theorem 1 states that with an egoistic motive ( $r > 0$ ), these cutoffs shift down. This means that there are belief regions in which individuals behave differently whether they have an egoistic motive or not. We illustrate an example in Figure 1. Precisely, when the current belief  $p_t$  is between the two lower cutoffs— $p_l(0)$  and  $p_l(r)$ —individuals continue acquiring information only with the egoistic motive. Similarly, only individuals with a egoistic motive stop acquiring information at beliefs between the two upper cutoffs  $p_h(r)$  and  $p_h(0)$ .

Figure 1: Stopping and continue intervals with and without egoistic motive



In an empirical population not all individuals may have the same underlying preference type. However, for all preference types, the *continue interval* of beliefs *below* the prior  $p_0$  is *larger* when the individual has an egoistic motive relative to when she does not. For a given belief  $p_t < p_0$ , this implies that the likelihood of the belief of a random individual being in the continue interval is *larger* when there is an egoistic motive.

Similarly, for all preference types, the *continue interval* of beliefs *above* the prior  $p_0$  is *smaller* when the individual has an egoistic motive relative to when she does not. For a given belief  $p_t > p_0$ , this implies that the likelihood of the belief of a random individual being in the continue interval is *smaller* when there is an egoistic motive.

We formally derive the following two predictions from Theorem 1 (see Appendix B).

**Prediction 1.** When individuals have an egoistic motive, conditional on having observed information so that  $p_t < p_0$ , they are more likely to *continue* acquiring further information, compared to when they do not have an egoistic motive.

**Prediction 2.** When individuals have an egoistic motive, conditional on having observed information so that  $p_t > p_0$ , they are more likely to *stop* acquiring further information, compared to when they do not have an egoistic motive.

As the Figure 1 illustrates, Theorem 1 does not necessarily predict *strictly* different behaviour for all beliefs: in the example of the figure, the given preference type continues for beliefs  $p_t$  close to the prior  $p_0$  whether having an egoistic motive or not. Empirically, this means: first, when we average differences of the stopping behaviour at all beliefs  $p_t > p_0$  or  $p_t < p_0$ , the theory predicts *strict* differences between the scenarios with and without egoistic motive. Second, when we compare the individuals' stopping decisions after having observed certain information—that is, at a given belief  $p_t$ —we might observe similar behaviour between the setting with an egoistic motive and the one without or strict differences. Notably, while similarities in the information choices would be consistent with the predictions, the predictions rule out observing the opposite of fishing for good news.

### 1.6.2 Fishing for good news may improve social welfare

When the agent has an egoistic motive, the model predicts that the agent alters her way of acquiring information (Prediction 1 and 2). Having different information in turn affects which option she chooses and the resulting harm on the other. We call this *indirect* effect of the egoistic motive on the welfare of the other the *information effect*.

How does the bias in information acquisition affect the welfare of the other? Our theory predicts that it may increase the other's welfare, i.e., the information effect may be positive.

**Prediction 3. (Information effect)**

There is an open set of agent types, so that the information effect of an egoistic motive  $r > 0$  is positive, i.e. it decreases the likelihood of making harmful decisions.

The information effect can only be positive if the agent does not acquire full information when there is no egoistic motive. Not acquiring full information without egoistic motive is optimal whenever  $l(a) \neq 1$ , which means that being more certain that the action  $a$  is harmless for the other does not improve the belief utility once the threshold level of certainty  $l(a)$  is reached. Prediction 3 captures that agents who are satisficers in this sense may be better informed at the decision-stage when they have an egoistic motive since this motive will make them fish for good news.

Besides the information effect, the egoistic motive might *directly* affect the agent's decision, given her belief. We call this the *decision effect*. We show that the decision effect is negative whenever the agent's belief utility only depends on the likelihood of harming the other.

**Prediction 4. (Decision effect)**

Take an agent with belief-based utility that only depends on the likelihood of harming the other. The decision effect of an egoistic motive  $r > 0$  is negative, i.e., it increases the likelihood of making harmful decisions.

A formal definition of the information effect and the decision effect are provided in Appendix B.6.

### 1.6.3 Fishing for good news as a strategic behaviour

We test a further prediction consistent with the strategic model of information acquisition.

**Prediction 5. (Intelligence)**

The tendency to fish for good news (Prediction 1 and 2) is *stronger* among more intelligent agents.

Prediction 5 is in line with our *strategic* model of fishing for good news. Evidence supporting it would suggest that fishing for good news is a strategic phenomenon.

#### 1.6.4 The predictions cannot be tested with static experiments with one-shot information

The predictions derived from the theory can *not* be tested with static experiments where the individuals can decide between receiving a given piece of information or no information.

In such experiments, the exogenous precision of the information determines the continuation payoff after the information choice, *ceteris paribus*. It, in turn, determines the decision makers' information choice. This poses two problems: first, the observed information choices will depend on the experiment design choices. Second, the *exogenous* continuation payoff is different from the equilibrium continuation payoff given optimal strategies. Hence, static experiments with one information decision cannot be used to investigate optimal information acquisition strategies.<sup>15</sup>

### 1.7 Discussion: how the belief-based utility can be interpreted

The key feature of the belief-based utility  $u$  in our model is that some beliefs are more desirable than others. Many psychological concepts share this feature: e.g., guilt aversion, self-image concerns or belief-based moral constraints, as we discuss below in detail. This paper's analysis of the prevalent motive to hold desirable beliefs hence highlights the *similarity* between different models and their prediction power, and we do not attempt to distinguish them in the later empirical analysis.

#### 1.7.1 Variation of the model with self-image concerns

We formulate a variation of the model in Section 1.1. In this variation, the belief-based utility  $u$  captures self-image concerns, a prominent explanation

---

<sup>15</sup>A concrete example: consider a one-shot experiment where individuals receive one piece of information with a relatively low precision. Suppose that utilities are concave in the belief. In this case, theory predicts that whenever prior beliefs are sufficiently low or sufficiently high, the piece of information will not affect the individuals' choice. Thus, the agent will avoid the piece of information since the belief utility is concave. However, given a design that avoids such censoring effects, theory would predict information seeking for low priors under the optimal strategy (see Theorem 1).

of willful ignorance (Grossman and van der Weele, 2017, see e.g.), and the belief-based utility arises in a signalling equilibrium (see Appendix D.2).

Bodner and Prelec (2003) propose the notion of *diagnostic utility*: an agent’s decisions are ‘diagnostic’ about her type. Building on this notion, Grossman and van der Weele (2017) provide a model with self-image concerns in prosocial decisions and a binary information choice. Our model variation follows this existing self-image concern literature in three ways: first, the agent has a continuous prosocial type  $\theta \in [0, 1]$  that captures how much she cares about the welfare of the other relative to her own remuneration. Second, she values her prosocial self-image, i.e., her utility depends on her belief about her prosocial type. Third, the final belief together with her action choice are diagnostic about her prosocial type.

We provide sufficient conditions for the existence of a *monotone* equilibrium (Theorem 2 in Appendix D.2): in this equilibrium, the least prosocial types avoid information about their action’s consequences altogether. The more prosocial types acquire some information about the state, and the higher their type, the strictly more information they acquire in a Blackwell sense.<sup>16</sup> Hence, agents self-signal high prosociality by stopping at more informative beliefs.

### 1.7.2 Other interpretations

**Moral constraints.** Rabin (1995) provides a model of moral constraints in which an agent has an exogenous constraint on pursuing her egoistic interest. Specifically, the agent maximizes her self-interest subject to the constraint that her action is not too likely to harm others, captured by a cutoff probability. When the belief-based utility  $u(a, q)$  of the model in Section 1.1 takes the following form (8), we may interpret it as capturing such internal moral constraints,<sup>17</sup>

$$u(a, q) = \begin{cases} -w & \text{if } q < l(a), \\ 0 & \text{if } q \geq l(a), \end{cases} \quad (8)$$

---

<sup>16</sup>A signal  $s$  about a binary state  $\omega \in \{\alpha, \beta\}$  is Blackwell more informative than a signal  $\tilde{s}$  if the distribution of posterior belief  $\Pr(\alpha|s)$  is a mean preserving spread of the distribution of the posterior  $\Pr(\alpha|\tilde{s})$ .

<sup>17</sup>Rabin (1994) considers the same type of utility function in his analysis.

where  $q$  is the agent’s belief about the likelihood that the action is harmless to the other and where  $w > r$ . The condition  $w > r$  ensures that the *moral constraint*  $q \geq l(a)$  is binding, meaning that the agent will not choose the self-benefiting option  $x$  unless her final belief  $q$  satisfies  $q \geq l(a)$ . The main result Theorem 1 holds for the specification (8).

**Guilt (aversion).** The belief-based utility  $u(a, q)$  from the model in Section 1.1 may capture emotions of guilt. Typically, guilt is formulated as a relative notion in games. That is, the agent’s guilt increases in the harm that she inflicts on the receiver, relative to some expectation of the receiver (see e.g., Battigalli and Dufwenberg, 2007). Loosely following the literature, a guilt formulation of the belief-based utility is<sup>18</sup>

$$u(a, q) = \begin{cases} q - l & \text{if } q < l, \\ 0 & \text{if } q \geq l, \end{cases} \quad (9)$$

where  $l \leq 1$  is the ‘receiver expectation’ of not being harmed and  $q$  is the agent’s belief about the likelihood that the agent’s choice  $a$  is harmless to the receiver. The main result Theorem 1 holds for the specification (9).

## 2 A laboratory experiment

We conduct a laboratory experiment with modified binary dictator games. All participants have the same initial endowment. Contingent on an unknown state, one of the two options that the dictator has to choose from reduces the payoff of the receiver, while the other option does not reduce the payoff of the receiver. Before deciding, the dictator can acquire costless information about which option has a negative externality on the receiver.

### 2.1 The treatment variations

The key treatment variation in our experiment is whether one option in the dictator game generates more payoff for the dictator than the other. In the ‘*Tradeoff*’ treatment, one option increases the dictator’s payoffs, while the

---

<sup>18</sup>See e.g., chapter 3.1 in Battigalli and Dufwenberg (2020).

other does not. In the ‘*Control*’ treatment, neither option affects the dictator’s payoffs. The comparison between *Tradeoff* and *Control* pins down the causal effect of having a self-benefiting option on the dictator’s information acquisition behaviour. We describe the details of this treatment variation below when we present the dictator game.<sup>19</sup>

## 2.2 The dictator game

	Good state ( $x$ harmless)	Bad state ( $y$ harmless)		Good state ( $x$ harmless)	Bad state ( $y$ harmless)
$x$	( <b>0</b> , 0)	( <b>0</b> , -80)	$x$	(+ <b>25</b> , 0)	(+ <b>25</b> , -80)
$y$	(0, -80)	(0, 0)	$y$	(0, -80)	(0, 0)

(a) *Control* Treatments

(b) *Tradeoff* Treatments

These tables present the dictator games in the *Control* and *Tradeoff* treatments. The number pairs in the table present (dictator’s payment, receiver’s payment).

Table 1: Dictator Decision Payment Schemes

At the beginning of the experiment, all subjects receive 100 experimental points as an endowment. Each experimental point is equivalent to 5 Eurocents. With this endowment, the subjects play the dictator game. Table 1 presents the payment scheme of the dictator game in *Tradeoff* and *Control* respectively. In both treatments, the dictator chooses between two options,  $x$  and  $y$ . There are two states of the world, ‘ $x$  harmless’ or ‘ $y$  harmless’. Depending on the state, either option  $x$  or option  $y$  reduces a receiver’s payment by 80 points, while the respective other option does not affect the receivers’ payment. Note that each option harms the receiver in one of the states. This design makes sure that the dictator cannot avoid the risk of harming the receiver without information about the state. In *Control*, the dictator receives no additional points regardless of her choice and the state. In *Tradeoff*,  $x$  is self-benefiting for the dictator: she receives 25 additional points when choosing  $x$ , but no additional points when choosing  $y$ .

<sup>19</sup>We implement a second treatment variation to address potential self-selection, which we detail in Section 4.2.

***Good state vs Bad state.*** For the ease of exposition, we hereafter refer to state ‘*x harmless*’ as the ‘*Good state*’, and state ‘*y harmless*’ as the ‘*Bad state*’. In state *x harmless*, option *x* increases in *Tradeoff* the dictator’s payments without disadvantaging the receiver. Believing that she is in state *x*, the dictator can choose the self-benefiting option *x* without feeling immoral. In contrast, in state *y harmless*, option *x* increases the dictator’s payment at the cost of a reduction in the receiver’s payment. Although this labeling is not meaningful in the *Control* treatments, we will generally refer to ‘*x harmless*’ as the *Good state* and ‘*y harmless*’ as the *Bad state* for consistency.

The dictator starts the experiment without knowing the state that she is in. She only knows that in every twenty dictators, seven are in the *Good state*, and thirteen are in the *Bad state*. That is, the dictator starts the experiment with a prior belief of 35% on that she is in the *Good state* and of 65% on that she is in the *Bad state*. A high prior belief in the *Bad State* strengthens the moral dilemma: choosing the self-benefiting option *x* without further information most likely harms the receiver. The prior belief is the same in *Control* and *Tradeoff*. Hence, the comparison between *Tradeoff* and *Control* is not driven by the asymmetric prior belief.

Before making the decision, the dictator can draw additional information and obtain more accurate beliefs about the state that she is in. We describe the information in the next subsection.

### 2.3 The noisy information

We design a noisy information generator for each state, which generates information that is easily interpretable. Specifically, each piece of information is a draw from a computerized box containing 100 balls. In the *Good state*, 60 of the balls are white and 40 are black; in the *Bad state*, 40 balls are white and 60 are black (see Figure 3 in Appendix A.1). The draws are with replacement from the box that matches each dictator’s actual state. After each draw, we display the Bayesian posterior belief about the likelihood of each state on the dictator’s individual computer screen, to reduce the cognitive cost of interpreting the information and to prevent non-Bayesian updating.

**Good news vs. bad news** For the ease of exposition, we refer to a white ball as a piece of ‘good news’ and a black ball as a piece of ‘bad news’. This is because, in the *Good state*, dictators draw a white ball with a higher probability. Hence, the draw of a white ball leads to an increase in the dictator’s belief about the likelihood of the *Good state*—the state in which the dictator in the *Tradeoff* treatments can choose  $x$  and gain the additional payment without reducing the payment of the receiver. Reversely, in the *Bad state*, dictators would draw a black ball with higher probability. A black ball is an evidence for the *Bad state*. In *Control*, we will still refer to a white ball as good news and a black ball as bad news for consistency, although the dictators in *Control* should not have a preference over the two states, hence also not over the color of the balls.

## 2.4 The experimental procedure

The experiment consists of three parts: the preparation stage, the main stage, and the supplementary stage.

**The preparation stage.** (i) The dictator reads paper-based instructions about the dictator decision and the noisy information. (ii) In these instructions, we also describe Bayes’ rule and tell the dictator that later in the experiment, we are going to help them to interpret the information by showing them the Bayesian posterior beliefs after each ball that they draw. (iii) Besides, the instructions specify that each experiment participant starts the experiment with 100 points of an endowment. (iv) We also inform them that option  $x$  is harmless for 7 out of 20 of the dictators and  $y$  for 13 out of 20. That is, the dictator’s prior information is that the *Good state* has a likelihood of 35% on the *Good state* and that the *Bad state* has a likelihood of 65%.

After reading the instructions, the dictators answer five control questions designed to check their understanding of the instructions. They keep the paper-based instructions for reference throughout the experiment.

**The main stage.** In the main stage, (i) the dictators can acquire information about the state that they are individually in; (ii) they choose between  $x$  and  $y$  in the dictator game.

Specifically, the dictator can acquire a piece of information by clicking a button that makes the computer draw a ball randomly from the box matched to their actual individual state (see Figure 3 in A.1). The draws are with replacement. After each draw, the screen displays the latest ball drawn and the Bayesian posterior beliefs about the *Good state* and the *Bad state* given all the balls drawn so far (rounded to the second decimal, see Figure 4 in Appendix A.1). There are two buttons on the screen: one to draw an additional ball, and the other to stop drawing and proceed to the dictator game. To proceed, the dictator must click on one of the buttons.

Besides drawing balls, throughout the experiment the dictators have no other way to learn about the true state that they are in. It is common knowledge that the receiver does not learn the information acquired by the dictator.

The draws do not impose any monetary cost on the dictator. The time cost of acquiring information is limited: between draws, there is a mere 0.3 second time lag to allow the ball and the Bayesian posterior belief to appear on the computer screen. It means that a dictator can acquire 100 balls within 30 seconds, which would yield almost certainty.

Having ended the information acquisition, the dictator chooses between  $x$  and  $y$  in the dictator game in Table 1a (in the *Control* treatments) or Table 1b (in the *Tradeoff* treatments). After that, in the implementation stage, the dictators' choices are implemented and the payments are calculated.

**The supplementary stage.** (i) We elicit the dictator's posterior beliefs about the state after the dictator game. The belief elicitation is incentivized.<sup>20</sup> We compare the elicited and the Bayesian posterior beliefs in Appendix C.4. We find that, for a majority of dictators, their elicited posterior beliefs and their Bayesian posterior beliefs coincide, and the deviation is not significantly different between *Tradeoff* and *Control* (two-sided Mann-Whitney-U test,  $p=0.29$ ). (ii) The subjects take part in the Social Value Orientation (SVO) slider test, which measures 'the magnitude of concern people have for others' and categorizes subjects into altruists, prosocials, individualists, and the competitive type (Murphy *et al.*, 2011). (iii) The subjects answer a questionnaire surveying their socio-demographics, e.g., the gender and age. They also answer a 5-item Raven's progressive matrices test (Raven *et al.*, 1998), which measures cognitive ability. We report the details of the supplementary stage

in Appendix C.4.

**Treatment assignment and implementation.** We randomize within each laboratory session: (i) the *Tradeoff* and *Control* treatments, (ii) the states: we randomly assign 35% of the laboratory terminals to the *Good state*, and 65% to the *Bad state*. The subjects are then randomly seated and randomly matched in a ring for the dictator game. The subjects are told that their decisions would affect the payment of a random participant in the same experimental session other than themselves. After all the subjects have made their dictator decisions, the experiment moves on to the implementation stage, where we inform the subjects that the dictator game decisions are being implemented and their payments are affected according to another participant’s dictator game decision. Each subject plays the dictator game only once.

We conducted the experiment in October and December 2018 at the BonnEconLab. 496 subjects took part (250 in *Tradeoff* and 246 in *Control*). Among the subjects, 60% are women, and 93% are students. They are, on average, 24 years old, with the youngest being 16 and the oldest being 69. The subjects are balanced between treatments, concerning gender, student status, and age (see Appendix C.4). We used z-tree (Fischbacher, 2007) to implement the experiment and hroot (Bock *et al.*, 2014) to invite subjects and to record their participation. Instructions and interfaces on the client computers were written in German, as all subjects were native German speakers.

**Payments.** In the experiment, payments are denoted in points. One point equals 0.05 Euro. At the end of the experiment, the details of the points and the equivalent payments earned in the experiment are displayed on the individual computer screens. The subjects received payments in cash before leaving the laboratory. The total earnings of a subject were the sum of the

---

<sup>20</sup>We incentivize the belief elicitation using the randomized Quadratic Scoring Rule adapted from Drerup *et al.* (2017) and Schlag *et al.* (2013). For the stated belief that the likelihood of the good state is  $b\%$ , we calculate the following value

$$M = \begin{cases} \frac{(b-100)^2}{100} & \text{if } x \text{ is harmless,} \\ \frac{b^2}{100} & \text{if } y \text{ is harmless.} \end{cases} \quad (10)$$

Then, the computer draws a random number  $A \sim U[0, 100]$  and the dictator receives 30 points if  $A > M$ .

following components: an endowment of 5 Euro, an additional 1.25 Euro if the subject was in treatments *Tradeoff* and chose  $x$ , a 4 Euro reduction if the subject’s randomly assigned dictator made a decision that reduces her payments, a random payment of either 1.5 or 0 Euro for revealing their posterior beliefs, a payment ranging from 1 to 2 Euro depending on the subject’s decisions in the SVO slider test, a payment ranging from 0.3 to 2 Euro depending on the decisions in the SVO slider measure of another random subject in the same laboratory session, and a fixed payment of 3 Euro for answering the questionnaire. A laboratory session lasted, on average, 45 minutes, with an average payment of 11.14 Euro.

### 3 Findings

In this section, we present the empirical analyses using data from our experiment. We provide summarizing statistics in Appendix C.1 and proceed below with the analyses of the dictators’ information acquisition behaviour. The median number of pieces of information acquired by the dictators is 5.

#### 3.1 Main finding: individuals fish for good news

The main finding from the experiment is that the dictators in *Tradeoff* ‘fish for good news’: Compared to the dictators in *Control*, having received more bad news, the dictators in *Tradeoff* are more likely to *continue* acquiring information; having received more good news, the dictators in *Tradeoff* are more likely to *stop* acquiring information.

In Section 3.1.1, we demonstrate *fishing for good news* by comparing between *Tradeoff* and *Control* the dictators’ decisions to continue acquiring information after the first draw and after the second draw. In Section 3.1.2, we analyze the entire information acquisition histories, leveraging statistical tools from survival analysis.

##### 3.1.1 Behaviour after the first pieces of information

To demonstrate the main result, we first consider the dictators’ decision to continue or stop acquiring information after having received one piece of in-

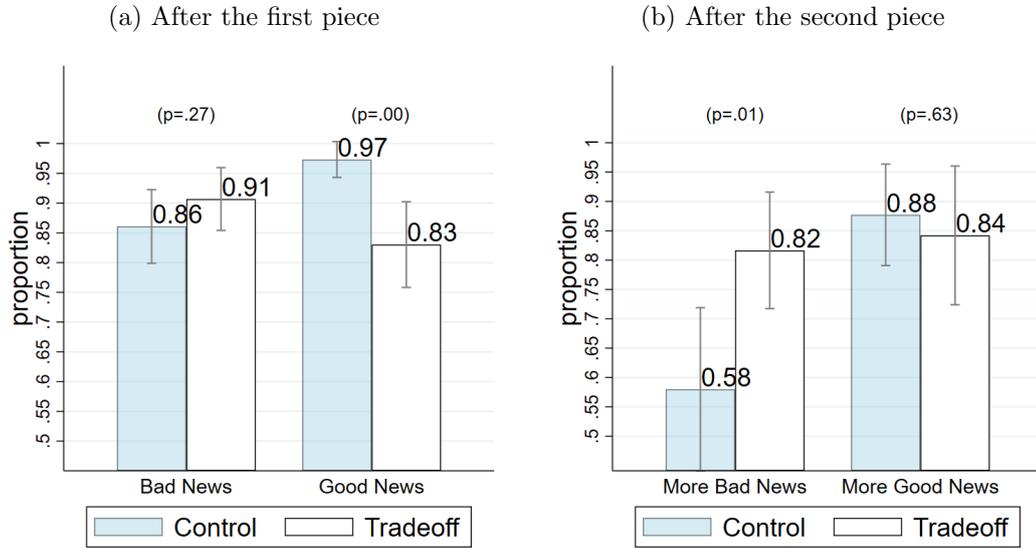
formation and two pieces of information.

**Finding 1** *The dictators' decision to continue acquiring information after the first piece of information differs from the decision in the Control baseline. Differences depends on the information that has been acquired. The same holds after the second piece of information. (i) Having received more bad news, weakly more dictators continue acquiring information in Tradeoff than in Control. (ii) Having receiving more good news, weakly less dictators in Tradeoff continue acquiring information than in Control.*

First, as predicted by *fishing for good news*, the treatment effect on whether to continue acquiring information depends on the information history. Specifically, in a logistic regression, the interaction effect of being in *Tradeoff* and having acquired more good news is significantly negative (after the first piece of information:  $p = .00$  and after the second piece of information:  $p = .049$ ; for details see Table 3 in Appendix A.2).

Specifically, the effects summarized in Finding 1 demonstrate '*fishing for good news*': compared to the *Control* baseline, facing a self-benefiting option makes individuals more likely to *continue* acquiring further information when the previous information indicates the that this option harms others. On the opposite, when the previous information suggests that the self-benefiting option is harmless, individuals are more likely to *stop* acquiring information. Figure 2 presents the proportions of dictators who continue acquiring information right after the first piece of information and the first two pieces of information, with the  $p$  values in the parentheses.

Figure 2: Proportion of dictators continuing after the first draws



This figure presents the proportion of dictators who continue acquiring information after the first piece of information (2a) and first two pieces of information (2b). Note that in *Control*, the within treatment difference given different first news is due to the asymmetric prior belief of 35% in the *Good* state. *Control* serves as the baseline to control for the effect of the prior belief on the information acquisition strategy.

Note that the theory predicts strict differences between treatments for some, but not for all given information sequences (a detailed discussion is made in Section 1.6.1). In particular, the similar effect after one piece of bad news is consistent with the theoretical predictions, so is the similar effect after two pieces of good news. Notably, the theory rules out the opposite of *fishing for good news* for all given information sequences, in line with the effects in Figure 2. Further, it predicts *strict* between-treatment comparisons aggregating over all information sequences. We will test this latter prediction in Section 3.1.2.

### 3.1.2 The entire information histories

We only provide simple comparative statics for the first two pieces of information since the sample size and the statistical power shrink as the information

process unfolds and some dictators stop acquiring information.<sup>21</sup> To jointly estimate the effect of having received more good news or more bad news aggregating over all information histories, we leverage tools from survival analysis instead.

**Model specification.** We carry out our analysis in the framework of the Cox proportional hazards model (Cox, 1972). The Cox proportional hazards model is often used for studying what influence individuals’ hazards of choosing an exit option when they face exit decisions repeatedly (e.g., in the literature on unemployment, see Card *et al.*, 2007; Michelacci and Ruffo, 2015, etc.). We use the Cox model to investigate the dictators’ hazards to stop acquiring information when they repeatedly decide whether to acquire further information. The Cox model has the advantage that it can address the dynamic selection that happens as observations drop out from the observed process. We will discuss this in detail in Section 4.2. Another advantage of the Cox model is that the coefficient estimates have a direct interpretation as hazards ratios, which we will explain momentarily when interpreting our estimation results.<sup>22</sup>

Taking *Control* as the baseline, we analyze the dependence of the *Tradeoff* dictators’ decision to continue or to stop acquiring information on whether up to that point they have received more good news or more bad news. The model specification is the following:

$$h(t|X) = h_0(t) \cdot \exp(\beta_1 \text{Tradeoff} + \beta_2 \text{Info} + \beta_{12} \text{Tradeoff} \times \text{Info} + \alpha z_t). \quad (11)$$

$h(t|X)$  denotes the dictator’s hazard rate to stop acquiring information after the  $t$ -th piece of information, given the set of covariates  $X$ ; the baseline hazard function  $h_0(t)$  captures the hazards over the draws at covariate vector 0.<sup>23</sup> The three covariates of interest are the treatment dummy “Tradeoff”; “Info”, the categorical variable denoting information histories that have more pieces of bad news, good news, or an equal number of bad and good news, with bad news dominance as the baseline; and the interaction of the two.

---

<sup>21</sup>In Section 4.2, we discuss how our empirical framework addresses potential issues with self-selection, explaining that our estimates are lower bounds for the effects.

<sup>22</sup>We report a robustness check using a logistic model in Appendix A.6.

<sup>23</sup>The Cox model naturally includes no constant term, since  $h_0(t)$  already captures the hazard rate at covariate vector 0 (see for example Cleves *et al.*, 2010).

**Model assumptions.** The model is correctly specified if (a) the covariates shift the baseline hazard proportionally, so that the hazard rate  $h(t|X)$  is multiplicative in the covariates (*proportional hazards assumption*), and (b) there are no omitted variables.

The proportional hazards assumption can be violated when some subgroups of the sample have different baseline hazards,  $h_0(t)$ . To make sure this assumption is not violated, stratification on the characteristics that might affect the hazard rate is often employed (see e.g., Blossfeld *et al.*, 2019). Stratification allows the baseline hazards  $h_0(t)$  to vary on the strata, while it estimates the aggregate effect of the covariates across all the stratified groups. We stratify our estimation on the following characteristics that can affect the baseline hazard: gender, cognitive ability (measured by the score in a Raven’s matrices test), and prosocial types (categorized by the SVO measure of Murphy *et al.*, 2011). After the stratification, the proportional hazards assumption of the Cox model is not violated, whereas without stratification it is.<sup>24</sup> <sup>25</sup> We also control for the accuracy of the individual belief after each ball drawn.<sup>26</sup>

It has been shown that omitting variables in the Cox model would only lead to underestimating the effects of interest (see Bretagnolle and Huber-Carol, 1988). To conclude, since the proportional hazards assumption is not violated in our estimation and since omitted variables can only lead to underestimation of the effects of interest, the significant results that we find are lower bounds of the effects.

**Data.** To test the dependency of the treatment effect on the information history, we first need to obtain a crucial independent variable: a factor variable denoting whether after a draw the information history has more good news or

---

<sup>24</sup>See Table 5 in Appendix A.4.

<sup>25</sup>We report a robustness check using a logistic model in Appendix A.6. The results of the logistic model are in line with those of the Cox model. The logistic model can be viewed as a hazards model with a proportional odds ratio assumption (Cox, 1975). However, unlike the Cox model, it does not allow for the baseline hazard to vary with the covariates. That is, it makes stronger assumptions than the stratified Cox model. Details are in the Appendix.

<sup>26</sup>In the experiment, the prior belief in the Good state is 0.35, a belief smaller than 0.5. Therefore, the posterior belief is more accurate after an information history with  $k$  more pieces of bad news than goods news relative to one with  $k$  more pieces of good news than bad news. To prevent this difference from being picked up by the Info dummy, we control for information accuracy. For this, we use the (expected) Brier score (Brier, 1950) of the beliefs as a proxy for the accuracy of beliefs:  $\text{belief}_{\text{Good}} \times \text{belief}_{\text{Bad}}^2 + \text{belief}_{\text{Bad}} \times \text{belief}_{\text{Good}}^2$ .

more bad news. Within an individual, this variable can vary after each draw. That is, the variable is time-varying. To obtain time-varying covariates, we follow the survival analysis literature and split each dictator’s information history at the unit of individual draws (see Blossfeld *et al.*, 2019, pp 137-152). The resulting data set consists of *pseudo-observations* at the person-draw level. For every draw of each dictator, the pseudo-observation records the dictator’s information history up to that draw and whether the dictator chooses to stop or continue acquiring information directly after that draw. For each pseudo-observation, we distinguish between information histories with more pieces of good news, more pieces of bad news, or the same number of good and bad news.

## Results.

**Finding 2** *Compared to the Control baseline, (i) having received more bad news than good news, the dictators in Tradeoff are more likely to **continue** acquiring information; (ii) while they are more likely to **stop**, having received more good news than bad news.*

We find that given information histories with more bad news than good news, being randomly assigned to *Tradeoff* has a significantly *negative* effect on the dictators’ hazards to stop acquiring information ( $\hat{\beta}_1 = -.29$ ,  $p = .02$ ). The interaction term between the treatment and having acquired more good news is significantly *positive* ( $\hat{\beta}_{12|good} = .43$ ,  $p = .03$ ). We proceed by explaining these results below and report the details of the Cox model estimation in Table 5 in Appendix A.4.

These findings show that the dictators fish for good news. First, given information histories with more bad news than good news, the between treatment comparison of the stopping hazard can be expressed by the following hazards ratio:

$$\begin{aligned}
 HR_{bad} &= \frac{h(t|bad, Tradeoff)}{h(t|bad, Control)} = \frac{\exp(\beta_1 \cdot 1 + \beta_2 \cdot 0 + \beta_{12} \cdot 1 \cdot 0 + \alpha_{z_t})}{\exp(\beta_1 \cdot 0 + \beta_2 \cdot 0 + \beta_{12} \cdot 0 \cdot 0 + \alpha_{z_t})} \\
 &= \frac{\exp(\beta_1 + \alpha_{z_t})}{\exp(\alpha_{z_t})} \\
 &= \exp(\beta_1); \tag{12}
 \end{aligned}$$

Recall that  $\hat{\beta}_1 = -.29 < 0$ . This means that when dictators have acquired more bad news, the hazard to stop is *lower* in *Tradeoff* than in *Control*. Specifically, the ratio between the two is  $\exp(-.29) \approx .75$ —in *Tradeoff* the hazard to stop acquiring information is 25% lower than in *Control*.

Second, given information histories with more good news than bad news, the between treatment comparison of the stopping hazard can be expressed by the following hazards ratio:

$$\begin{aligned} \text{HR}_{\text{good}} &= \frac{h(t|\text{good}, \text{Tradeoff})}{h(t|\text{good}, \text{Control})} = \frac{\exp(\beta_1 \cdot 1 + \beta_{2|\text{good}} \cdot 1 + \beta_{12|\text{good}} \cdot 1 \cdot 1 + \alpha z_t)}{\exp(\beta_1 \cdot 0 + \beta_{2|\text{good}} \cdot 1 + \beta_{12|\text{good}} \cdot 0 \cdot 1 + \alpha z_t)} \\ &= \frac{\exp(\beta_1 + \beta_{2|\text{good}} + \beta_{12|\text{good}} + \alpha z_t)}{\exp(\beta_{2|\text{good}} + \alpha z_t)} \\ &= \exp(\beta_1 + \beta_{12|\text{good}}). \end{aligned} \tag{13}$$

Since  $\hat{\beta}_1 = -.29$  and  $\hat{\beta}_{12} = .43$ ,  $\exp(\beta_1 + \beta_{12|\text{good}}) = \exp(-.29 + .43) \approx 1.15$ . That is, the hazard to stop acquiring information is 15% *larger* in *Tradeoff* than in *Control* when the dictators have received more good news than bad news up to that point.

## 3.2 Fishing for good news may improve receiver welfare

In this section, we analyse how an egoistic motive affects the welfare of the receivers. Notably, in our experiment, we implement a baseline treatment, *Control*, where the dictators have no egoistic motive in the decision. This allows us to identify causal effects of having an egoistic motive on social welfare.

In Section 3.2.1, we find that fishing for good news *improves* receiver welfare in our experiment (Finding 3). This finding stands in stark contrast to the existing empirical literature on information choices and motivated beliefs, which has observed that people *avoid* information; intuitively, information avoidance can only *harm* others. In Section 3.2.2, we provide some intuition for our result.

### 3.2.1 Analysis

In *Tradeoff*, option  $x$  increases the dictator's payment while option  $y$  does not. This egoistic motive causes the dictators to fish for good news (Finding

1 and 2). Does fishing for good news make the dictators in *Tradeoff* more often choose the option that reduces the receiver’s payment?

Directly comparing the receiver welfare between *Tradeoff* and *Control* confounds the two effects predicted by the model (see Section 1.6.2): first, for a given belief about the harmful option, the egoistic motive in *Tradeoff* may affect the dictator’s *decision* between the two options (*the decision effect*); second, the egoistic motive may affect her belief about the harmful option and in turn her decision, by making her fish for good news (*the information effect*).

Our theoretical analysis predicts that the decision effect *decreases* the likelihood of receivers’ being spared from harm; and in contrast, the information effect may *improve* receiver welfare (see Section 1.6.2 and Appendix B.6 for details). Below, we develop an empirical strategy to disentangle the decision and the information effect in our experimental data.

**Identification Strategy.** To empirically disentangle the two effects, we construct a *Counterfactual* scenario, in which the dictators acquire information as in *Control*, but decide as in *Tradeoff* given the final posterior beliefs (as illustrated in Table 2). When comparing the receiver welfare in the *Counterfactual* to the *Control* treatment, we isolate the decision effect by keeping fixed the final posterior beliefs; when comparing the receiver welfare in the *Counterfactual* to that in the *Tradeoff* treatment, we isolate the information effect by keeping fixed the decision between  $x$  and  $y$  given beliefs.

Table 2: Constructing the *Counterfactual* scenario

	<i>Tradeoff</i>	<i>Control</i>
posterior beliefs		×
decision given belief	×	
compared to the <i>Counterfactual</i>	<i>information effect</i>	<i>decision effect</i>

**Finding 3** *The information effect is positive: having a self-benefiting option changes how dictators acquire information in a way that they more often choose the harmless option, controlling for the dictators’ decision given their belief.*

First, we compare *Tradeoff* with the *Counterfactual* and find a positive information effect. In *Tradeoff*, the proportion of unharmed receivers is higher than in the *Counterfactual* (68% compared to 62%, Chi-Square  $p = .046$ ). This shows that had the dictators in *Tradeoff* acquired information the way the dictators in *Control* did, they would have inflicted *more harm* on the receivers.

**Finding 4** *The decision effect is negative: controlling for the dictators' beliefs, having a self-benefiting option makes the dictators less often choose the harmless option.*

Second, we compare the *Counterfactual* with the *Control* and find a negative decision effect. In the *Counterfactual*, the proportion of unharmed receivers is lower than in the *Control* treatment (62% compared to 73%, Chi-Square  $p = .00$ ). Consistent with this finding, our data show that in both states the dictators in *Tradeoff* are more likely to choose  $x$ , the self-benefiting option, than dictators in *Control* (Chi-Square,  $p = .00$ ).

Aggregating the two effects, the proportion of unharmed receivers does not significantly differ between *Tradeoff* and *Control* (68% compared to 73%, Chi-Square  $p = 0.17$ ).

### 3.2.2 Intuition

In *Control*, where the dictators' own payments are unaffected by their decisions, if the dictators would acquire full information they would not cause harm to the receivers.<sup>27</sup> However, the dictators in *Control* only acquire a limited amount of information, leaving room for fishing for good news to improve social welfare, as in Finding 3.<sup>28</sup>

In fact, in *Control*, 27% of the dictators choose the option that eventually harms the receivers. The dictators in *Control* only acquire a limited amount of information—the median number of pieces of information that they draw is 5.

---

<sup>27</sup>In our data, all the 7 *Control* dictators who acquire information until the Bayesian posterior beliefs displayed to them are rounded to certainty cause no harm to the dictator.

<sup>28</sup>In a recent paper by Exley and Kessler (2021), the authors make a similar observation: even when an egoistic motive is missing, a substantial fraction of dictators avoid information revealing the consequences of their options.

### 3.3 Dictators with higher IQ have a higher tendency to fish for good news

In Section 1, we formally show that the trade-off between belief utility and material rewards leads dictators in *Tradeoff* to *fish for good news*. In line with this strategic explanation, we find that the effects in Finding 2 are *larger* among the dictators with a cognitive ability above the median (measured by a Raven’s matrices IQ test), compared to among all the dictators.

**Finding 5** *The dictators with an above median IQ have a higher tendency to fish for good news.*

Specifically, we split the sample at the median of the dictators’ scores in a Raven’s matrices test and estimate the Cox model as in (11) using the respective data. When only the dictators with above-median IQ are considered, the estimates of our coefficients of interest,  $\beta_1$  and  $\beta_{12|good}$ , are significant at the 1% and 2% level. In particular, having received more bad news, the above-median dictators’ hazards to stop acquiring information in *Tradeoff* is 32% *lower* than that in *Control*—an effect size that is 28% larger compared to the effect when all dictators are considered.<sup>29</sup> Having received more good news, the hazard to stop acquiring information in *Tradeoff* is 27% *higher* than that in *Control*—an effect that is 80% larger than the effect when all dictators are considered.<sup>30</sup> In contrast, when we only consider the dictators with below-median IQ, the estimated coefficients of interest,  $\beta_1$  and  $\beta_{12|good}$ , are both insignificant. We report the details of these analyses in Table 6 in Appendix A.5.

## 4 Discussion

### 4.1 Bayesian persuasion

In this section, we discuss the relation of our model in Section 1.1 to the literature on Bayesian persuasion (Kamenica and Gentzkow, 2011). The equi-

---

<sup>29</sup>The hazard ratio between the *Tradeoff* and *Control*,  $HR_{bad} = \exp(\beta_1)$ , is estimated to be  $\exp(-.38) = .68$ .

<sup>30</sup>The hazard ratio between the *Tradeoff* and *Control*,  $HR_{Good} = \exp(\beta_1 + \beta_{12,Good})$ , is estimated to be  $\exp(-.38 + .62) = 1.27$ .

librium characterization in Section 1.2 (Lemma 1) shows that any equilibrium can be characterized by a Bayes-consistent distribution of stopped beliefs  $p_\tau$  with support on two posteriors  $p_l \leq p_0 \leq p_h$ . This distribution maximizes

$$E(V(p)) \tag{14}$$

for  $V(p) = \max_{a \in \{x,y\}} U(a, p; r)$  across all Bayes-consistent distributions of posterior beliefs.

This formulation of the (dynamic) information acquisition problem of the agent with two competing motives makes the relation to models of *interpersonal* Bayesian persuasion most apparent. The problem is equivalent to that of a sender who tries to ‘persuade’ a distinct other agent by transmitting information about a payoff relevant state to her and where the sender’s payoff only depends on the posterior belief  $p$  of the other agent and is given by  $V(p)$  (compare to Kamenica and Gentzkow, 2011). Hence, on a conceptual level, the problem of using information to align two internal motives, and the problem of transmitting information to another person with the goal to align the person’s beliefs and actions with the sender’s interest, are analogous. We note that Schweizer and Szech (2018) apply techniques from Bayesian persuasion, to study the optimal revelation of medical information in a setup with anticipatory utility.

**Persuasion cutoffs in the data.** We impute the empirical distribution of the cutoffs  $p_h$  and  $p_l$  from our experiment data. Figure 5 in Appendix A.7 shows the empirical distributions. We find that among the responsive types, the prediction of Theorem 1 holds, that is, the distributions of the cutoff  $p_h$  and  $p_l$  are shifted downwards when there is an egoistic motive (one-sided Kolmogorov-Smirnov test,  $p = .045$  for  $p_h$  and  $p = 0.074$  for  $p_l$ ).

## 4.2 Self-selection

Below we explain the self-selection facing our empirical analysis and show evidence that self-selection can only *weaken* our finding of fishing for good news. That is, our results are lower bounds for the effects. We also provide intuition based on the model for why this is the case.

In our experiment, individuals repeatedly decide whether to continue or to stop acquiring information. When one compares between treatments the stopping decision at any point of the information process, dictators who have already stopped earlier are not in the sample. That is, the sample is dynamically self-selected. This could confound the empirical analysis if the self-selected samples differ between treatments.

To directly address this potential confound, we implement a second treatment variation: *Force* and *NoForce*. While in ‘*NoForce*’, the dictators are *not forced* to acquire any information, in *Force* the information stage starts for all dictators, and the dictators acquire at least one piece of information. Therefore, the sample for our analysis of the stopping decision after the first piece of information is *not* self-selected in *Force*. Similar to Finding 1, in this sample, we also find a significantly negative interaction effect between *Tradeoff* treatment and having acquired more good news on the dictators tendency to continue acquiring information (logistic regression  $p = .05$ , see Table 4 in Appendix A.3 for details).

Besides, to estimate the effects across all information histories, we use the Cox proportional hazard model. When there are omitted variables in the Cox model, dynamic selection may confound the analysis; however, Bretagnolle and Huber-Carol (1988) show that this can only lead to *underestimating* the effects of interest. In our setup, this result is intuitive, as we outline below with an argument based on our theory.

To illustrate why the self-selection weakens the observed “fishing for good news” behaviour, let’s first consider the selection into the information process in the *NoForce* treatment variation. Here, in *Tradeoff*, 25 out of 26 dictators who do not acquire information choose  $x$ , while in *Control*, 10 out of 12 dictators who do not acquire information choose  $y$ . Had these dictators received a further piece of information *supporting* their dictator decisions, i.e., good news in *Tradeoff* and bad news in *Control*, theory predicts that they would also stop directly and take the same decisions, an effect that would strengthen Finding 1. This prediction follows directly from the cutoff structure of optimal strategies, which are given by beliefs  $p_l \leq p_h$ .<sup>31</sup> This (theoretical) argument

---

<sup>31</sup>For example, if an agent prefers to stop at the prior and choose  $x$ , she would also stop and choose  $x$  at any higher belief  $p_t > p_0 = p_h$ . An agent who would stop at the prior and choose  $y$ , would also stop and choose  $y$  at any lower belief  $p_t < p_0 = p_l$ .

extends to the decisions about further pieces of information, suggesting that the potential self-selection could only weaken the results about the entire information histories (Finding 2).

In summary, the potential issue of dynamic self-selection could only have weakened our result instead of having driven it. Further, in a sample clean of self-selection, we find similar results as in Finding 1.

### 4.3 Information cost

In this paper, we examine how an individual's desire for believing in her good moral conduct affects her way of acquiring relevant information. The focus of our investigation is the tradeoff between an egoistic motive and a morally desirable belief. In emphasizing individuals having preferences over beliefs, our investigation follows the motivated belief literature (Bénabou and Tirole, 2016). In this section, we discuss another factor that may affect information acquisition, information cost.

Information cost might take different forms, e.g., material cost, cognitive cost or time cost of sampling. Methodologically, to avoid our empirical results being confounded by information cost, we are careful to limit the information cost in our experiment. First, the subjects do not pay for the information in any material form. Second, we limit the subjects' cognitive costs of interpreting information by giving them the Bayesian posterior beliefs after each draw of information. If cognitive costs drove our results, the results should be weaker for those with higher cognitive ability (see Bénabou and Tirole, 2016). In our experiment, Finding 5 shows that when restricting the analysis to subjects with above-median cognitive ability, our main result (Finding 2) becomes *stronger*, a result showing that our main finding is unlikely driven by cognitive costs. Third, we minimize the time cost of sampling by imposing only a minimum time lag of 0.3 second between draws of information. On average, the dictators in the experiment spend 57 seconds acquiring information.

While our investigation focuses on the trade-off between beliefs and egoistic rewards, the role of cost for information acquisition is an interesting topic studied, for example, in the literature of rational inattention. The interaction between beliefs, rewards and information costs might be an interesting

direction for future research.

## 5 Concluding remarks

Theoretically and experimentally, this paper analyzes the effect of the trade-off between a motive to feel moral and a competing egoistic motive on individuals' information acquisition strategies.

Two features of our study stand out relative to the existing literature. First, we consider environments with rich information sources. This ensures that the predicted and the observed information choices are not confounded by exogenous limitations imposed by the study. Theoretically, it means that we characterize the globally optimal information acquisition strategy. Empirically, it allows us to observe uncensored data on the individuals' information strategies and uncover novel phenomena. Second, we consider a baseline in which the egoistic motive is removed from the decision. Comparing this baseline with the scenario with an egoistic motive, we can study the causal effects of the trade-off between competing motives.

Three main findings emerge: first, having competing motives makes individuals *fish for good news*. Specifically, individuals are more likely to *continue* acquiring information after receiving information suggesting negative externalities of a selfish decision ('bad news'). Reversely, after receiving information indicating no harm ('good news'), individuals are more likely to stop. Second, theoretically, *fishing for good news* may improve social welfare. This prediction is supported by our data. Finally, the tendency to *fish for good news* is stronger among individuals with above-median cognitive ability—evidence that fishing for good news is more likely a strategic behaviour than a result of cognitive limitations.

The paper opens up directions of future research. The key feature of our setting is that the decision-maker has two competing motives—one urging the individual to choose a particular action, and the other urging her to act upon her belief about an unknown state. Such a trade-off is present in many economic contexts beyond moral decisions. Imagine a food lover presented with a delicious new dish. While she longs for the dish, she also wants to believe that the food that she consumes is healthy. How would she inquire about the

nutrition facts of the food? Similar trade-offs arise for example in smoking and workout decisions.<sup>32</sup> Besides, it might be interesting to investigate information acquisition strategies in field settings where morality and egoism might clash. Last, the theoretical model can be used to study other questions about individuals' *information preferences*—a recently active area of empirical research.<sup>33</sup> For example, one may analyze preferences over information skewness in settings where individuals have two competing motives.

---

<sup>32</sup>We note that theoretical work by Carrillo and Mariotti (2000) considers a setup where two self-related motives arise from time inconsistent preferences and show that strategic ignorance may result. Schweizer and Szech (2018) consider a setup with a medical patient whose utility depends both on her future health outcomes and her anticipation. They show that the patient-optimal medical test partially reveals her health condition.

<sup>33</sup>See e.g. Masatlioglu *et al.* (2017) and Falk and Zimmermann (2016).

# Appendix

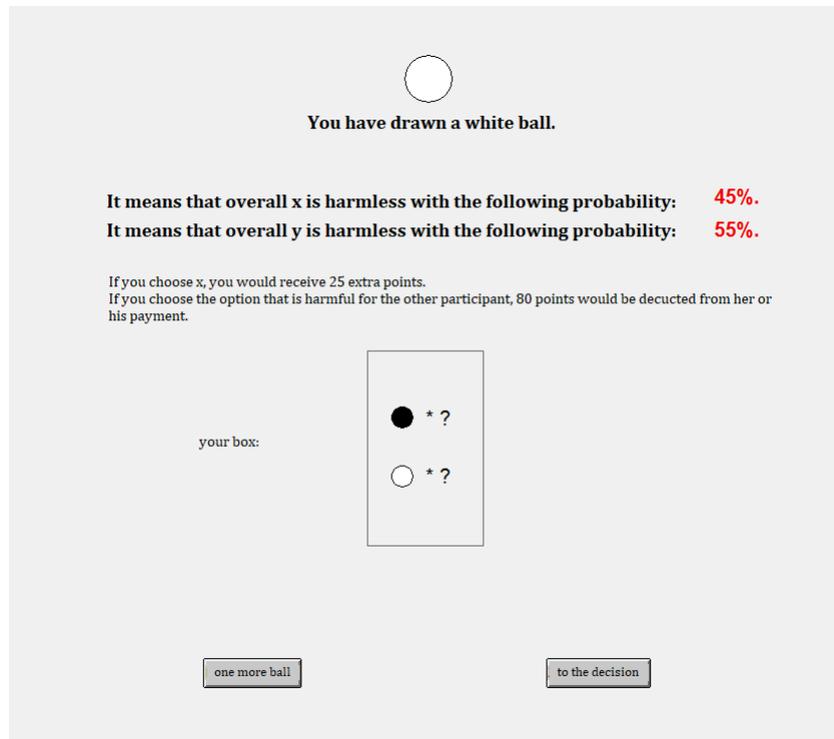
## A Empirical appendices

### A.1 Experimental design

Figure 3: The Noisy Information Generators



Figure 4: Screenshot of the Information Stage



## A.2 The difference in difference in Finding 1

To test the difference in the difference in Finding 1, we estimate the following logistic regression for the probability to continue acquiring information after the first piece of information and the first two pieces of information:

$$\text{logit}(\text{continue}) = b_1 \text{Tradeoff} + b_2 \text{good} + b_{12} \text{Tradeoff} \cdot \text{Good} + c, \quad (15)$$

where Tradeoff is a factor variable for the treatment; good is a factor variable for whether the dictator has acquired more good news or bad news. Table 3 presents the regression estimates.

The interaction effect between *Tradeoff* and having acquired more good news, i.e.,  $b_{12|\text{good}}$ , is significantly negative, showing that the treatment effect on the probability to continue acquiring information significantly differs after having received more good and after having received more bad news.

Table 3: Logistic regression estimates (with  $p$  in the brackets)

Coef.	(1)	(2)
	After the 1st piece of info	After the first 2 pieces of info
$b_1$	.45 (.27)	.117 (.008)
$b_{2 \text{good}}$	1.77 (.006)	1.64 (.001)
$\mathbf{b_{12 \text{good}}}$	<b>-2.46</b> <b>(.001)</b>	<b>-1.46</b> <b>(.049)</b>
$c$	1.82 (.00)	.32 (.26)
N	458	409
Chi2 p	.00	.00

### A.3 Difference in difference in *Force*

Below we estimate the logistic regression in (15) in the *Force* treatment, where there is no self-selection in the sample. We find that consistent with Finding 1, having receiving more good news and being randomly assigned to *Tradeoff* has a significantly negative interaction effect on the dictators' tendency to continue acquiring information.

Table 4: Logistic regression estimates (with  $p$  in the brackets)

Coef.	(1)
$b_1$	.48 (.53)
$b_{2 good}$	1.52 (.17)
$\mathbf{b_{12 good}}$	<b>-2.74</b> <b>(.039)</b>
$c$	2.03 (.00)
N	161
Chi2 p	.05

## A.4 The Cox model results

Table 5: The Cox proportional hazards model results (with  $p$  in brackets)

Coef.	Covariate	(1)	(2)
$\hat{\beta}_1$	<b>Tradeoff</b>	-.29 (.02)	-.21 (.09)
$\hat{\beta}_{12}$	<b>Tradeoff</b> $\times$ <b>Good news dominance</b>	.43 (.03)	.28 (.14)
	Balanced	-.35 (.35)	-.53 (.15)
$\hat{\beta}_2$	Good news dominance	-.14 (.38)	-.09 (.59)
	Balanced	-.52 (.03)	-.51 (.03)
Stratified by: gender, IQ, prosociality		Yes	No
<b>Violation of the proportional hazards assumption</b>		<b>No</b>	<b>Yes</b>
Control variable: belief accuracy		Yes	Yes
Observations (individuals)		458	458
Chi2 p-value		.00	.00

This table presents the estimated *coefficients* of the Cox model in (11), with standard errors clustered at the individual level. In the brackets, we report the p value of the corresponding coefficient estimate. The dependent variable is the hazard to stop acquiring information, and the key coefficients of interests are  $\hat{\beta}_1$  and  $\hat{\beta}_{12}$ .  $\exp(\hat{\beta}_1)$  reflects the treatment effect on the dictators' hazards to stop acquiring further information, given information histories dominated by bad news; and  $\exp(\hat{\beta}_1 + \hat{\beta}_{12}|\text{Good news dominance})$  reflects the treatment effect on the hazards, given information histories dominated by good news (see the derivation in Equation (13)). The violation of the proportional hazard assumption of the Cox model (PH) is tested using Schoenfeld residuals. Without stratification, the PH assumption is violated, as shown in column (2), implying that the baseline hazard might differ for subgroups of the sample. Hence, we follow the literature and use stratification to allow the baseline hazard to vary according to the control variables, i.e., gender, the prosocial types (categorized by the SVO test), and the cognitive ability (categorized by the score in a 5-element Raven's matrices test). With the stratification, PH is no longer violated. We also control for the belief accuracy, measured by the Brier score of the beliefs after each draw (see Footnote 26). The reported likelihood Chi-square statistic is calculated by comparing the deviance ( $-2 \times \log\text{-likelihood}$ ) of each model specification against the model with all covariates dropped. We use the Breslow method to handle ties.

## A.5 Individuals with higher IQ have a higher tendency to fish for good news

Table 6: The Cox model results for about and below median IQ (with  $p$  in brackets)

Coef. Covariate	Above Median (1)	Below Median (2)	All Dictators (3)
$\hat{\beta}_1$ <i>Tradeoff</i>	-.38 (.01)	-.22 (.25)	-.29 (.02)
$\hat{\beta}_{12}$ <b>Tradeoff</b> $\times$ <b>Good news dominance</b>	.62 (.02)	.24 (.43)	.43 (.03)
Balanced	.31 (.55)	-1.00 (.08)	-.35 (.35)
$\hat{\beta}_2$ Good news dominance	-.14 (.51)	-.25 (.31)	-.14 (.38)
Balanced	-1.01 (.01)	-.23 (.47)	-.52 (.03)
Stratified	Yes	Yes	Yes
<b>Violation of the PH assumption</b>	No	No	No
Control variable: belief accuracy	Yes	Yes	Yes
Observations (individuals)	267	191	458
Chi2 p-value	.00	.00	.00

This table presents the Cox model results for the dictators with above and below median cognitive ability, measured by the number of correctly answered questions in a Raven’s matrices test. Standard errors are clustered at the individual level. For comparison, we include the result for the whole sample in Column (3). In Column (1) and (2), the estimation is stratified by gender and prosociality, but not by IQ since we explicitly compare the dictators with above and below median IQ here. In Column (3), the estimation is stratified by gender, prosociality and IQ. The median number of correct answers to the Raven’s test is four out of five in our experiment. In this table, the dictators above the median have given correct answers to four or five questions in the Raven’s test, and the subjects below the median have correctly answered less than four questions in the Raven’s test. The finding is that dictators with a higher cognitive ability have a higher tendency to “fish for good news”. For a comprehensive description of the Cox model estimation, please see the description of Table 5.

## A.6 Robustness check of the Cox model estimate: a logistic regression

Using the data at the person-draw level, we estimate the following logistic model as a robustness check of the Cox model estimate and find a result

similar to Finding 2 from Section 3.1.2.

$$\text{logit } h(X) = X_t \cdot b + Z \cdot a + (C + T \cdot c), \quad (16)$$

where  $h(X)$  is the probability that the dictator stops acquiring information after that draw;  $X$  denotes the same covariates of interest as in the Cox model, i.e.,

$$X \cdot b = \beta_1 \text{Tradeoff} + \beta_2 \text{Info} + \beta_{12} \text{Tradeoff} \times \text{Info}. \quad (17)$$

The control variables in  $Z$  include gender, cognitive ability, prosociality and belief accuracy, all measured in the same way as in the Cox model in Section 3.1.2.  $T$  is a vector of time dummies, which captures the time dependency of the probability to stop acquiring information.

When interpreting the results, this logistic model can be viewed as a hazard model in which the covariates proportionally affect the *odds* of stopping the information acquisition (Cox, 1975). Formally, consider such a hazard model,

$$\frac{h(t)}{1 - h(t)} = \frac{h_0(t)}{1 - h_0(t)} \cdot \exp(X_t \cdot b + Z \cdot a).$$

Then,

$$\underbrace{\log\left(\frac{h(t)}{1 - h(t)}\right)}_{\text{logit } h(X)} = \underbrace{\log\left(\frac{h_0(t)}{1 - h_0(t)}\right)}_{C+T \cdot c} + X_t \cdot b + Z \cdot a. \quad (18)$$

Unlike in the framework of the Cox model, the coefficients here cannot be interpreted as hazard ratios. Instead, they should be interpreted as odds ratios. Our prediction that the hazard to stop acquiring information is lower in *Tradeoff* when bad news dominates suggests a negative  $\beta_1$  in (17). And the prediction that the hazard is higher when good news dominates suggests a positive  $\beta_1 + \beta_{12|good}$ . Results reported in Table 7 support these predictions.

Table 7: The logistic model results (with  $p$  in the brackets)

Coef. Covariate	(1)	(2)
$\hat{\beta}_1$ <b>Tradeoff</b>	-.26 (.08)	-.22 (.126)
$\hat{\beta}_{12}$ <b>Tradeoff</b> $\times$ <b>Good news dominance</b>	.36 (.099)	.35 (.08)
Balanced	-.54 (.18)	-.59 (.14)
$\hat{\beta}_2$ Good news dominance	-.19 (.30)	-.62 (.00)
Balanced	-.66 (.02)	-1.11 (.00)
Control	Yes	No
N (person-draws)	4,658	4,658
Pseudo R2	.07	.05

This table presents the estimated coefficients of the logistic model in (16), with standard errors clustered at the individual level. The dependent variable the decision to stop acquiring information, and the key coefficients of interests are  $\hat{\beta}_1$  and  $\hat{\beta}_{12}$ .  $\exp(\hat{\beta}_1)$  reflects the treatment effect on the dictator’s odds to stop acquiring further information, given information histories dominated by bad news. And  $\exp(\hat{\beta}_1 + \hat{\beta}_{12}|\text{good})$  reflects the treatment effect on the odds, given information histories dominated by good news. We control for belief accuracy, gender, the prosocial types (categorized by the SVO test), and the cognitive ability (measured by a Raven’s matrices test). The time dependency of the odds is accounted for by including a dummy for each period.

## A.7 Imputed Cutoff Distributions

In this section, we infer from the data the distribution of the optimal strategies, characterized by upper and lower belief cutoffs. In particular, we compare them between treatments.

Recall that in the model, we analyze the information acquisition behaviour of an agent who does not avoid information completely. Theorem 1 considers ‘responsive’ types, i.e., those who choose  $x$  if they stop at a posterior weakly above the prior or  $y$  if they stop at a posterior weakly below the prior. In our data, considering the subjects who do not avoid information completely, we find that the large majority of subjects behaves responsively (405 out of 458; *Control*: 225 out of 234; *Tradeoff*: 180 out of 224).<sup>34</sup>

<sup>34</sup>In the *Control* treatment, 5 dictators choose  $y$  after having received more good news,

Figure 5 shows the empirical cumulative distribution functions of the lower belief cutoff  $p_l$  (5b) and the upper belief cutoff  $p_h$  (5a). Both CDFs reflect the dictators who acquire some information. The CDF of the upper belief cutoff reflects the stopped beliefs of the dictators who stop weakly above the prior and choose  $x$ . The CDF of the lower belief cutoff reflects the stopped beliefs of the dictators who stop information acquisition at posterior beliefs weakly below the prior and choose  $y$ .

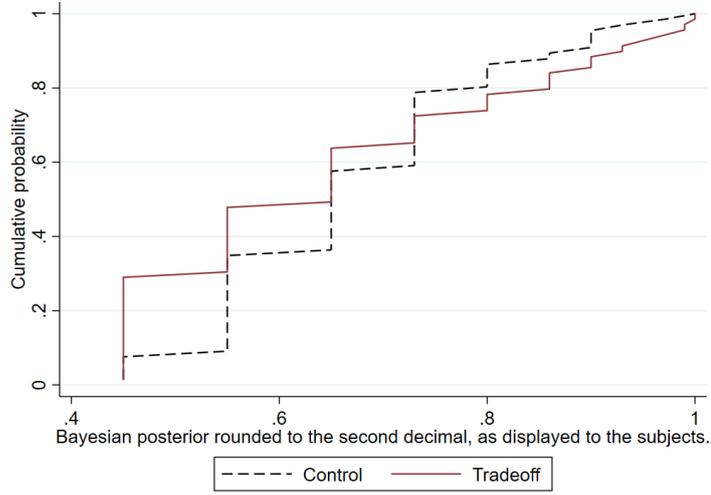
The figures show that the belief cutoffs are systematically lower in *Tradeoff*, as predicted by the model in Theorem 1 (one-sided Kolmogorov-Smirnov test,  $p = .045$  for  $p_h$  and  $p = 0.074$  for  $p_l$ ).

---

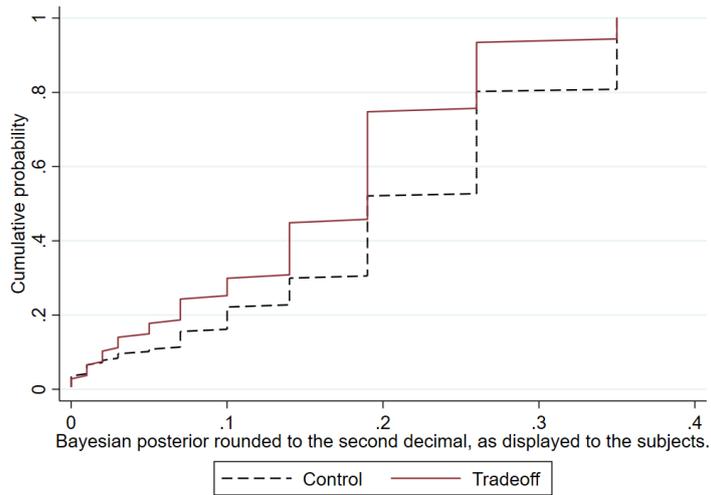
4 dictators choose  $x$  after having received more bad news. In the *Tradeoff* treatment, 4 dictators choose  $y$  after having received more good news, 37 subjects choose  $x$  after having received more bad news.

Figure 5: The CDF of the imputed cutoffs

(a) Upper cutoffs  $p_h$



(b) Lower cutoffs  $p_l$



## B Theory appendices

### B.1 Preliminaries for the proofs

First, we establish two claims that we will use to prove both Lemma 1 and Lemma 2. For this, recall the definition of the cutoff beliefs  $p_l$  and  $p_h$  following

the statement of Lemma 2.

**Claim 1** *Let  $p_t \in [p_l, p_h]$ . For any continuation strategy  $\tau$ ,*

$$\mathbb{E}(V(p_\tau)|(Z_s)_{s \leq t}) \leq \bar{V}(p_t) \quad (19)$$

**Proof.** We have

$$\mathbb{E}(V(p_\tau)|(Z_s)_{s \leq t}) \leq \mathbb{E}(\bar{V}(p_\tau)|(Z_s)_{s \leq t}) \leq \bar{V}(\mathbb{E}(p_\tau)|(Z_s)_{s \leq t}) = \bar{V}(p_t),$$

where we used that  $V \leq \bar{V}$  for the first inequality and Jensen's inequality for the second inequality. For the final equality, we use that  $\mathbb{E}(p_\tau|(Z_s)_{s \leq t}) = p_t$  by Doob's optional stopping theorem.<sup>35</sup> ■

Now, consider the candidate equilibrium strategy  $\tau^*$  where the agent continues to observe the information process as long as  $p_l < p_t < p_h$ , and stops whenever  $p_t \leq p_l$  or  $p_t \geq p_h$ .

**Claim 2** *Let  $p_t \in [p_l, p_h]$ . The strategy  $\tau^*$  satisfies*

$$\mathbb{E}(V(p_{\tau^*})|(Z_s)_{s \leq t}) = \bar{V}(p_t) \quad (20)$$

**Proof.** We consider two cases: if  $V(p_0) = \bar{V}(p_0)$ , by definition,  $p_h = p_l = p_0$  and the agent immediately stops at  $t = 0$ , i.e.,  $\Pr(p_{\tau^*} = p_0) = 1$ , which directly yields the result in this case. If  $V(p_0) < \bar{V}(p_0)$ , then,  $\bar{V}$  is linear on all open intervals  $I' \subseteq [\epsilon, 1 - \epsilon]$  satisfying  $p_0 \in I'$  and  $V(p) < \bar{V}(p)$  for all  $p \in I'$ , by its minimality. Now,  $(p_l, p_h)$  is the largest such interval, which implies that  $V$  and  $\bar{V}$  must coincide at  $p_l$  and  $p_h$ ,<sup>36</sup>

$$V(p_h) = \bar{V}(p_h), \quad \text{and} \quad V(p_l) = \bar{V}(p_l). \quad (21)$$

Finally, for any  $p_t \in [p_l, p_h]$ ,

$$\begin{aligned} \mathbb{E}(V(p_{\tau^*})|(Z_s)_{s \leq t}) &= \Pr(p_{\tau^*} = p_h|(Z_s)_{s \leq t})V(p_h) + \Pr(p_{\tau^*} = p_l|(Z_s)_{s \leq t})V(p_l) \\ &= \Pr(p_{\tau^*} = p_h|(Z_s)_{s \leq t})\bar{V}(p_h) + \Pr(p_{\tau^*} = p_l|(Z_s)_{s \leq t})\bar{V}(p_l) \\ &= \bar{V}(p_t), \end{aligned}$$

<sup>35</sup>See e.g., Revuz and Yor (2013).

<sup>36</sup>One checks that this is also true if  $(p_l, p_h) = (\epsilon, 1 - \epsilon)$  by the minimality of  $\bar{V}$ .

where we used (21) for the equality on the second line. For the equality on the third, we used the earlier observation that  $\bar{V}$  is linear on  $(p_l, p_h)$  together with Bayes' law. ■

## B.2 Proof of Lemma 1

Let  $\tau^*$  be the candidate equilibrium strategy where the agent continues to observe the information process as long as  $p_l < p_t < p_h$ , and stops whenever  $p_t \leq p_l$  or  $p_t \geq p_h$ . Claim 1 and Claim 2 together imply that at any point of time, following  $\tau^*$  is weakly optimal, hence  $\tau^*$  is an equilibrium. This proves Lemma 1.

## B.3 Proof of Lemma 2.

Take any equilibrium  $\tau^{**}$  in which the agent stops observing the information process whenever he is indifferent between stopping and continuing. It follows from Claim 1 and Claim 2 that, when  $p_t \in (p_l, p_h)$ , it is strictly optimal for the agent to continue acquiring information: stopping yields  $V(p_t)$ , which is strictly smaller than  $\bar{V}(p_t)$ , and there is a continuation strategy which yields  $\bar{V}(p_t)$  by Claim 2. When  $p_t \in \{p_l, p_h\}$ , if the agent would stop acquiring information, his payoff would be  $V(p_t) = \bar{V}(p_t)$ , given (21). Thus, it follows from Claim 1 that it is weakly optimal to stop acquiring information, so the agent stops under  $\tau^{**}$ . Finally, we conclude that  $\tau^{**}$  is identical to  $\tau^*$  (see the proof of Lemma 1 for the definition of the equilibrium  $\tau^*$ .)

## B.4 Proof of Lemma 3

Take the strategy  $\tau'$  where the agent never stops observing the information process (unless  $p_t \leq \epsilon$  or  $p_t \geq 1 - \epsilon$ , and she has to stop). Given  $\epsilon \approx 0$ , she acquires almost complete information about the state. Note that her expected utility when doing so is  $E(V(p_{\tau'})) \approx (1 - p_0)V(0) + p_0V(1) \geq (1 - p_0)u(y, 1) + p_0(u(x, 1) + r)$  since she can almost always choose  $y$  in the state when  $y$  is harmless and  $x$  in the state when  $x$  is harmless. Given that  $u(x, 1) = 0$  and  $u(y, 1) \leq 0$ , we have  $(1 - p_0)V(0) + p_0V(1) \geq (1 - p_0)u(y, 1) + p_0r > u(y, 1)$ . It follows that the equilibrium strategy  $\tau^*$  given by the cutoff beliefs  $p_l$  and

$p_h$  must yield a payoff strictly larger than  $u(y, 1)$  as well when  $\epsilon$  is sufficiently small, that is  $E(V(p_{\tau^*})) > u(y, 1)$ .

First, this implies that the agent does not choose  $y$  at  $p_h$  when  $\epsilon$  is sufficiently small: suppose she does so, then, she will also choose  $y$  at  $p_l < p_h$  since at  $p_l$  she is more certain that  $y$  is harmless, *ceteris paribus*. However, when she always chooses  $y$ , her payoff is weakly smaller than  $u(y, 1)$  since  $U(y, p, r) = u(y, 1 - p) \leq u(y, 1)$  for all  $p$ .

Second, this implies that  $V(p_h) > u(y, 1)$  when  $\epsilon$  is sufficiently small: suppose that  $V(p_h) \leq u(y, 1)$ . Then, also  $V(p_l) = \max_{a \in \{x, y\}} U(a, p_l, r) \leq u(y, 1)$  since  $U(y, p, r) = u(y, 1 - p) \leq u(y, 1)$  for all  $p$  and  $U(x, p_l, r) \leq U(x, p_h, r) \leq u(y, 1)$ . However,  $V(p_h) \leq u(y, 1)$  and  $V(p_l) \leq u(y, 1)$  together imply  $E(V(p_{\tau^*})) \leq u(y, 1)$ , which contradicts with the observation  $E(V(p_{\tau^*})) > u(y, 1)$  when  $\epsilon$  is small enough. Given that we assumed that the agent weakly prefers  $y$  at  $p_l$ , we have  $V(p_l) = u(y, 1 - p_l) \leq u(y, 1)$ . We conclude that  $V(p_h) > V(p_l)$  since  $V(p_h) > u(y, 1)$ .

## B.5 Derivation of Prediction 1 and 2

We derive two testable predictions from Theorem 1. These predictions hold whenever a large part of the population is of a ‘responsive’ type.<sup>37</sup>

To see why Prediction 1 holds, take an agent whose belief at time  $t > 0$  satisfies  $p_t < p_0$ . Let us compare the belief cutoff  $p_l$  of a randomly drawn ‘responsive’ preference type of this agent when  $r > 0$ , relative to when  $r = 0$ . Recall that Theorem 1 states that, for all responsive types,  $p_l$  is weakly smaller when  $r > 0$ . Hence, it is more likely that a randomly drawn type has the cutoff  $p_l < p_t$  instead of  $p_l \geq p_t$ . This makes it more likely that a random type continues to acquire information at  $t > 0$ . In a similar vein, Theorem 1 implies our second prediction.

---

<sup>37</sup>A responsive type uses information in a ‘responsive’ way, i.e., chooses  $y$  after information indicating that  $y$  is harmless to the other, and  $x$  after information indicating that  $x$  is harmless to the other. In our data 409 out of 459 subjects are responsive.

## B.6 Receiver welfare

### B.6.1 Definition of the information and the decision effect

Let  $v(a, \omega)$  denote the utility of the other when the agent chooses  $a \in \{x, y\}$  in  $\omega \in \{X, Y\}$ . For any  $p \in (0, 1)$  and  $r > 0$ , let  $a(p, r) = \arg \max_{a \in \{x, y\}} U(a, p, r)$ .<sup>38</sup> For any  $r > 0$ , let  $\tau(r)$  be the equilibrium information acquisition strategy of the agent given by the belief cutoffs  $p_l(r)$  and  $p_h(r)$ ; see Lemma 2 and thereafter, here we highlight the dependence on  $r$ . Given this notation,

$$\text{DE} = \mathbb{E} \left[ v(a(p_{\tau(0)}, r), \omega) \right] - \mathbb{E} \left[ v(a(p_{\tau(0)}, 0), \omega) \right]. \quad (22)$$

is the decision effect of the remuneration  $r > 0$  on the welfare of the other, and

$$\text{IE} = \mathbb{E} \left[ v(a(p_{\tau(r)}, r), \omega) \right] - \mathbb{E} \left[ v(a(p_{\tau(0)}, r), \omega) \right]. \quad (23)$$

is the information effect.

### B.6.2 Derivation of Prediction 3 and Prediction 4

Proposition 1 show that the decision effect is always negative when the belief-based utility only depends on the likelihood of harming the other (Prediction 4). It also shows that for an open set of types the information effect is positive (Prediction 3) and even offsets the decision effect, thereby leading to an overall positive effect on the welfare of the other. The proof is in Section B.6.3 and B.6.4.

**Proposition 1** *Take any  $r > 0$ .*

1. *For any preference type  $u$  with  $u(x, -) = u(y, -)$ , the decision effect is negative,  $\text{DE} \leq 0$ .*
2. *There is an open set of types  $u$ , so that the information effect and the overall effect are positive,  $\text{IE} > 0$  and  $\text{DE} + \text{IE} > 0$ .*

---

<sup>38</sup>When the agent is indifferent between  $x$  and  $y$ , let  $a(p, r) = x$ . The agent never stops when being indifferent, so that this choice is irrelevant for our analysis.

### B.6.3 Proof of Proposition 1, first item

Let  $p^*(r)$  solve

$$u(x, p) + r = u(y, 1 - p). \quad (24)$$

Let  $u(x, p) = u(y, p)$ . It is easy to see that  $p^*(r) < 0.5$  when  $r > 0$  and  $p^*(r) = 0.5$  when  $r = 0$ . Note that the agent's optimal decision between  $x$  and  $y$  is given by

$$a(p, r) = \begin{cases} x & \text{if } p > p^*(r) \\ y & \text{if } p < p^*(r). \end{cases} \quad (25)$$

Fix  $\bar{r} > 0$ . Comparing the scenarios when  $r = \bar{r}$  and when  $r = 0$ , the agent may only take different decisions if her final belief  $p_\tau$  is in  $[p^*(\bar{r}), 0.5)$ . At a final belief  $p_\tau \in [p^*(\bar{r}), 0.5)$ , the agent chooses  $x$  when  $r = \bar{r}$  and  $y$  when  $r = 0$ . The receiver is (weakly) better off with the choice  $y$  since the likelihood of harming the receiver with  $y$  is  $p_\tau < 0.5$  and the likelihood of harming the receiver with  $x$  is  $1 - p_\tau > 0.5$ . This shows that, fixing any information acquisition strategy  $\tau$ , the receiver is better off when the agent decides according to the decision rule  $a(-, 0)$  as opposed to the rule  $a(-, \bar{r})$ .

To prove the second item of Proposition 1, we take the preference types given by

$$u(x, q) = u(y, q) = \begin{cases} 0 & \text{for } q \geq l, \\ -\alpha(l - q) & \text{for } q < l \end{cases} \quad (26)$$

for some  $\frac{1}{2} < l < 1$  satisfying  $p_0 \in (1 - l, l)$ . We establish a preparatory result.

**Claim 3** *For all  $r \geq 0$ , it holds  $p^*(r) \rightarrow \frac{1}{2}$  when  $\alpha \rightarrow \infty$ .*

**Proof.** Note that  $\alpha \rightarrow \infty$  implies that, for  $q < l$ , we have  $u'(a, q) = \alpha \rightarrow \infty$ . In particular, for any  $\epsilon > 0$  there is  $\bar{\alpha}(\epsilon)$  such that for all  $\alpha > \bar{\alpha}$  and  $p < \frac{1}{2} - \epsilon$ , we have  $u(x, 1 - p) - u(x, p) > r$ . Hence,  $u(y, 1 - p) - u(x, p) > r$ , using that  $u(x, q) = u(y, q)$  for all  $q$ . Given (24), this implies  $p^*(r) \geq \frac{1}{2} - \epsilon$ . Further, it follows from (26) and (24) that  $p^*(r) \leq \frac{1}{2}$ , finishing the proof of the claim. ■

### B.6.4 Proof of Proposition 1, second item

Take a preference type as in (26). Recall the characterization of  $p_l(r)$  and  $p_h(r)$  after Lemma 1, where we highlight the dependence on  $r$  here. Take  $\bar{r} > 0$ . When  $\epsilon \approx 0$ ,

$$p_l(\bar{r}) = \epsilon, \text{ and} \quad (27)$$

$$p_h(\bar{r}) = l, \quad (28)$$

given (26), and Lemma 5 implies

$$p_l(0) = 1 - l, \text{ and} \quad (29)$$

$$p_h(0) = l \quad (30)$$

for  $\epsilon$  sufficiently small. First, it follows from Claim 3 that when  $\alpha$  is sufficiently large, both when  $r = 0$  and  $r = \bar{r}$ , the agent strictly prefers  $y$  at the belief  $p_l(0) = 1 - l < \frac{1}{2}$ . In other words, the decision rules  $a(-, \bar{r})$  and  $a(-, 0)$  both choose  $y$  at  $p_l(0)$ . Given (29) and the assumption that  $l > \frac{1}{2}$ , we have  $p_h(0) > \frac{1}{2}$ . Hence,  $p_h(0) > p^*(r)$  for all  $r \geq 0$ . So, the decision rules  $a(-, \bar{r})$  and  $a(-, 0)$  both choose  $x$  at  $p_h(0)$ . We conclude that the decision effect is zero.

Second, given (27) - (30), it holds  $p_h(\bar{r}) = p_h(0)$  and  $p_l(\bar{r}) < p_l(0)$  when  $\epsilon \approx 0$ . So, the agent chooses the option  $y$  at a higher belief in the harmfulness of  $y$  when  $r = \bar{r}$  relative to when  $r = 0$ . Further, the agent chooses the option  $x$  at the same belief. Fixing the decision rule  $a(-, 0)$  and varying the information acquisition strategy, we see that the information effect is strictly positive.

Altogether, we have shown for the types as in (26) that the sum of the decision effect and the information effect is strictly positive when  $\alpha$  is sufficiently large,  $\text{DE} + \text{IE} > 0$ . Clearly, by continuity, this is still true for any type with belief-based utility function  $\hat{u}$  ‘close-by’ to the belief-based utility  $u$  as in (26).<sup>39</sup>

---

<sup>39</sup>Formally, if we consider the  $L_2[0, 1]$ -norm on the space of belief-based utility functions, there is an open set  $U$  with  $u \in U$ , so that the sum of the decision effect and the information effect is positive for all  $\hat{u} \in U$ .

## References

- AKERLOF, G. A. and KRANTON, R. E. (2000). Economics and identity. *The Quarterly Journal of Economics*, **115** (3), 715–753.
- BATTIGALLI, P. and DUFWENBERG, M. (2007). Guilt in games. *American Economic Review*, **97** (2), 170–176.
- and — (2020). *Belief-dependent motivations and psychological game theory*. Working paper.
- BÉNABOU, R. and TIROLE, J. (2006). Incentives and prosocial behavior. *American Economic Review*, **96** (5), 1652–1678.
- and TIROLE, J. (2016). Mindful economics: The production, consumption, and value of beliefs. *Journal of Economic Perspectives*, **30** (3), 141–64.
- BLOSSFELD, H.-P., ROHWER, G. and SCHNEIDER, T. (2019). *Event History Analysis with Stata*. Routledge.
- BOCK, O., BAETGE, I. and NICKLISCH, A. (2014). Hroot: Hamburg registration and organization online tool. *European Economic Review*, **71**, 117–120.
- BODNER, R. and PRELEC, D. (2003). Self-signaling and diagnostic utility in everyday decision making. *The Psychology of Economic Decisions*, **1**, 105–26.
- BRETAGNOLLE, J. and HUBER-CAROL, C. (1988). Effects of omitting covariates in cox’s model for survival data. *Scandinavian Journal of Statistics*, pp. 125–138.
- BRIER, G. W. (1950). Verification of forecasts expressed in terms of probability. *Monthly weather review*, **78** (1), 1–3.
- BURDETT, K. (1996). Truncated means and variances. *Economics Letters*, **52** (3), 263–267.
- CARD, D., CHETTY, R. and WEBER, A. (2007). The spike at benefit exhaustion: Leaving the unemployment system or starting a new job? *American Economic Review*, **97** (2), 113–118.
- CARRILLO, J. D. and MARIOTTI, T. (2000). Strategic ignorance as a self-disciplining device. *The Review of Economic Studies*, **67** (3), 529–544.
- CLEVES, M., GOULD, W., GUTIERREZ, R. and MARCHENKO, Y. (2010). *An Introduction to Survival Analysis Using Stata*. College Station, TX, Stata Press.

- COX, D. R. (1972). Regression models and life-tables. *Journal of the Royal Statistical Society: Series B (Methodological)*, **34** (2), 187–202.
- (1975). Partial likelihood. *Biometrika*, **62** (2), 269–276.
- DANA, J., WEBER, R. A. and KUANG, J. X. (2007). Exploiting moral wiggle room: Experiments demonstrating an illusory preference for fairness. *Economic Theory*, **33** (1), 67–80.
- DI TELLA, R., PEREZ-TRUGLIA, R., BABINO, A. and SIGMAN, M. (2015). Conveniently upset: Avoiding altruism by distorting beliefs about others’ altruism. *American Economic Review*, **105** (11), 3416–42.
- DITTO, P. H. and LOPEZ, D. F. (1992). Motivated skepticism: Use of differential decision criteria for preferred and non-preferred conclusions. *Journal of Personality and Social Psychology*, **63** (4), 568.
- DRERUP, T., ENKE, B. and VON GAUDECKER, H.-M. (2017). The precision of subjective data and the explanatory power of economic models. *Journal of Econometrics*, **200** (2), 378–389.
- EIL, D. and RAO, J. M. (2011). The good news-bad news effect: Asymmetric processing of objective information about yourself. *American Economic Journal: Microeconomics*, **3** (2), 114–38.
- EXLEY, C. and KESSLER, J. B. (2018). *Motivated Errors*. Working paper.
- EXLEY, C. L. and KESSLER, J. B. (2021). *Information Avoidance and Image Concerns*. Tech. rep., National Bureau of Economic Research.
- FALK, A. and SZECH, N. (2016). *Pleasures of Skill and Moral Conduct*. Working paper.
- and ZIMMERMANN, F. (2016). *Beliefs and Utility: Experimental Evidence on Preferences for Information*. Working paper.
- FEILER, L. (2014). Testing models of information avoidance with binary choice dictator games. *Journal of Economic Psychology*, **45**, 253–267.
- FISCHBACHER, U. (2007). z-tree: Zurich toolbox for ready-made economic experiments. *Experimental economics*, **10** (2), 171–178.
- GEANAKOPOLOS, J., PEARCE, D. and STACCHETTI, E. (1989). Psychological games and sequential rationality. *Games and Economic Behavior*, **1** (1), 60–79.
- GINO, F., NORTON, M. I. and WEBER, R. A. (2016). Motivated bayesians: Feeling moral while acting egoistically. *Journal of Economic Perspectives*, **30** (3), 189–212.

- GNEEZY, U., SACCARDO, S., SERRA-GARCIA, M. and VAN VELDHUIZEN, R. (2016). *Motivated Self-Deception, Identity and Unethical Behavior*. Working paper.
- GOLMAN, R., HAGMANN, D. and LOEWENSTEIN, G. (2017). Information avoidance. *Journal of Economic Literature*, **55** (1), 96–135.
- GROSSMAN, Z. (2014). Strategic ignorance and the robustness of social preferences. *Management Science*, **60** (11), 2659–2665.
- and VAN DER WEELE, J. (2017). Self-image and willful ignorance in social decisions. *Journal of the European Economic Association*, **15** (1), 173–217.
- HAGENBACH, J. and KOESSLER, F. (2021). Selective memory of a psychological agent.
- HAISLEY, E. C. and WEBER, R. A. (2010). Self-serving interpretations of ambiguity in other-regarding behavior. *Games and Economic Behavior*, **68** (2), 614–625.
- KAMENICA, E. and GENTZKOW, M. (2011). Bayesian persuasion. *American Economic Review*, **101** (6), 2590–2615.
- KÖSZEGI, B. (2006). Ego utility, overconfidence, and task choice. *Journal of the European Economic Association*, **4** (4), 673–707.
- KUNDA, Z. (1990). The case for motivated reasoning. *Psychological bulletin*, **108** (3), 480.
- LOEWENSTEIN, G. (1987). Anticipation and the valuation of delayed consumption. *The Economic Journal*, **97** (387), 666–684.
- MASATLIOGLU, Y., ORHUN, A. Y. and RAYMOND, C. (2017). *Preferences for Non-Instrumental Information and Skewness*. Working paper.
- MICHELACCI, C. and RUFFO, H. (2015). Optimal life cycle unemployment insurance. *American Economic Review*, **105** (2), 816–59.
- MOBIUS, M. M., NIEDERLE, M., NIEHAUS, P. and ROSENBLAT, T. S. (2011). *Managing Self-Confidence: Theory and Experimental Evidence*. Tech. rep., National Bureau of Economic Research.
- MURPHY, R. O., ACKERMANN, K. A. and HANDGRAAF, M. (2011). Measuring social value orientation. *Judgment and Decision Making*, **6** (2), 771–781.
- RABIN, M. (1994). Cognitive dissonance and social change. *Journal of Economic Behavior & Organization*, **23** (2), 177–194.

- (1995). *Moral Preferences, Moral Constraints, and Self-Serving Biases*. Working paper.
- RAVEN, J. C. *et al.* (1998). *Raven's Progressive Matrices and Vocabulary Scales*. Oxford psychologists Press.
- REVUZ, D. and YOR, M. (2013). *Continuous Martingales and Brownian Motion*, vol. 293. Springer Science & Business Media.
- SCHLAG, K. H. *et al.* (2013). Eliciting probabilities, means, medians, variances and covariances without assuming risk neutrality. *Theoretical Economics Letters*, **3** (1), 38–42.
- SCHWEIZER, N. and SZECH, N. (2018). Optimal revelation of life-changing information. *Management Science*, **64** (11), 5250–5262.
- SERRA-GARCIA, M. and SZECH, N. (2019). *The (in) Elasticity of Moral Ignorance*. Working paper.
- SIMON, H. A. (1955). A behavioral model of rational choice. *The Quarterly Journal of Economics*, **69** (1), 99–118.
- ZIMMERMANN, F. (2020). The dynamics of motivated beliefs. *American Economic Review*, **110** (2), 337–61.

# Online appendix

## C Empirical appendices

### C.1 Summarizing statistics

Table 8: Basic information of the dictators

	no. obs.	Good State	female	student	av. age
<i>Tradeoff</i>	250	.35	.59	.94	24
<i>Control</i>	246	.36	.61	.93	24
p-value		.82	.62	.56	.49

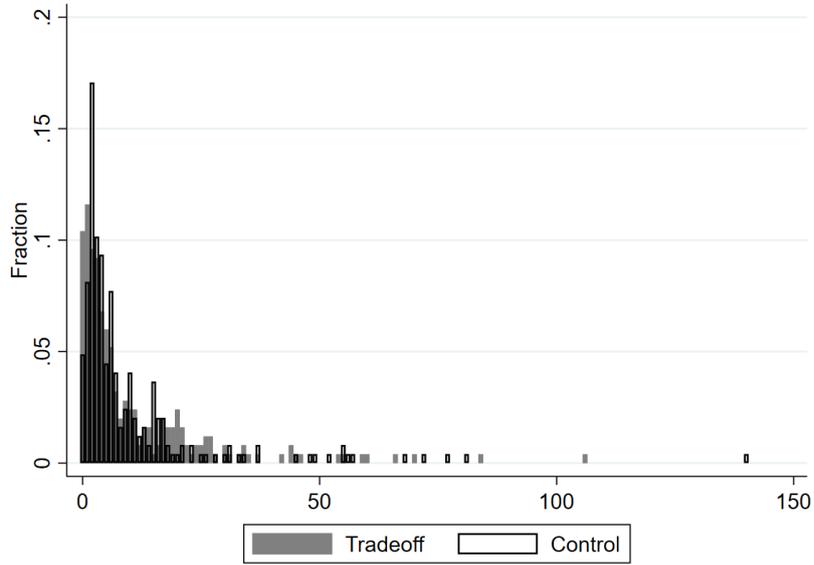
This table summarizes the basic characteristics of the dictators in each treatment. We compare these characteristics between *Tradeoff* and *Control*. For the state, gender, and student status we report the p-values of the Chi-Square test. For the dictators' age, we report the p-value of the two-sided t test.

Table 9: Information acquisition behaviour

	median no. balls	av. belief at decision
<i>Tradeoff</i>	6	.33
<i>Control</i>	5	.34
p-value	.24	.30

This table presents the median number of information pieces drawn by the dictators and their average Bayesian posterior beliefs. The p-values are of the two-sided Mann-Whitney-U test comparing between *Tradeoff* and *Control*.

Figure 6: Histogram of the number of information pieces drawn



This figure presents the number of information pieces that the dictators drew in *Tradeoff* and *Control* respectively. The distribution is not significantly different between *Tradeoff* and *Control* (two-sided Mann-Whitney-U test,  $p = .98$ ).

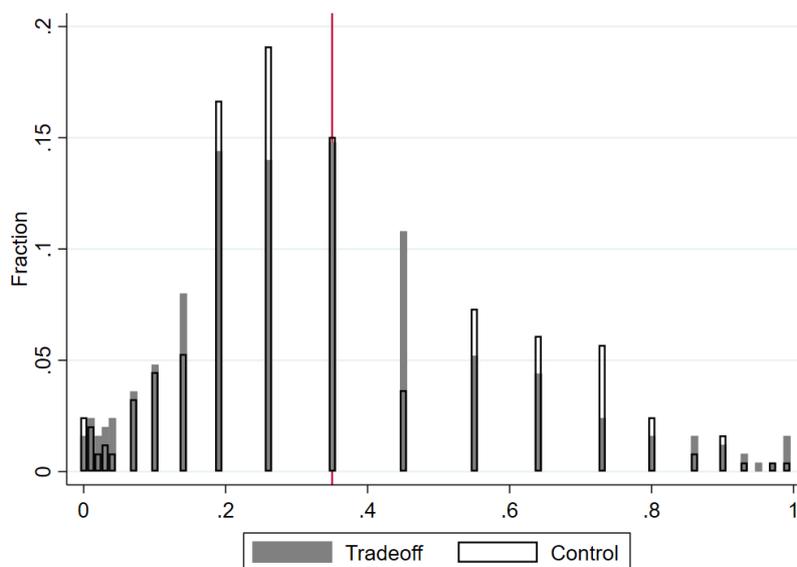
## C.2 The posterior beliefs

In this section, we investigate the dictators' final Bayesian posterior beliefs in the likelihood of the remunerative option being harmless. We compare the belief distributions between *Tradeoff* and *Control*. Our analysis is motivated by the main empirical results (Finding 1 and Finding 2) and Theorem 1. We find that having drawn more good news than bad news, dictators in *Tradeoff* stop acquiring information at a lower belief in the Good state than those in *Control* (0.63 in *Tradeoff*, 0.66 in *Control*, Mann-Whitney-U test  $p = 0.10$ ). In contrast, having drawn more bad news than good news, dictators in *Tradeoff* stop acquiring information at a similar belief in the Bad state than those in *Control* (0.84 in *Tradeoff*, 0.84 in *Control*, Mann-Whitney-U test  $p = 0.13$ ). Figure 7 presents the histogram of the dictators' posterior beliefs.

This asymmetry in the posterior beliefs mirrors the asymmetry in the information acquisition strategy in Finding 2. It is also in line with Theorem 1.

When it comes to the overall posterior beliefs, Bayes' consistency states that the mean posterior beliefs must be equal to the prior belief in both treatments. In accordance with it, we find that the mean posterior beliefs are not significantly different between *Tradeoff* and *Control* (mean: 0.33 in *Tradeoff*, 0.34 in *Control*; student t-test,  $p = 0.56$ ).<sup>40</sup>

Figure 7: Histogram of the posterior beliefs



This figure shows the histogram of the dictators' Bayesian posterior beliefs about the likelihood of the Good State when they end the information stage in *Tradeoff* and *Control*. The red vertical line represents the prior belief.

<sup>40</sup>This finding is reconciled with the distributional differences by the observation that in *Tradeoff* slightly more dictators ended up with a posterior belief above the prior, although the difference is insignificant (*Tradeoff*: 30%; *Control*: 29%; Chi-Square,  $p = .14$ ).

### C.3 Dictator Game Decision

Table 10: Dictator game decisions

	Choosing $x\%$			Harm %
	Good	Bad	Overall	
<i>Tradeoff</i>	.82	.40	.54	.32
<i>Control</i>	.55	.16	.30	.27
p-value	.00	.00	.00	.17

The first three columns of this table present the proportions of dictators who choose  $x$  in *Good* and *Bad* state and in each treatments. Recall that in the *Good* state,  $x$  does not harm the receiver, while in the *Bad* state it does. The last column presents the percentage of dictators whose decision reduce the receivers' payoffs in the dictator game. The p-values are from the Chi square tests comparing between *Tradeoff* and *Control* respectively.

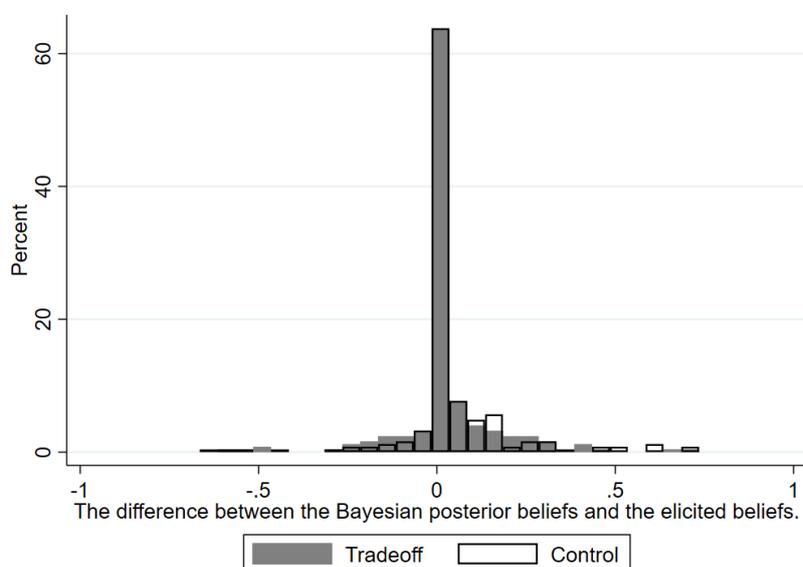
### C.4 The supplementary stage

After the experiment, we elicited the dictators' posterior beliefs on the state and their SVO scores. We also asked them to answer a questionnaire consisting of questions on their sociodemographics (gender, age, etc.). Five items from Raven's progressive matrices intelligence test are also included.

**Elicited beliefs** In the experiment, we display to the dictators the Bayesian posterior belief on the states (rounded to the second decimal) after each draw of information. After a dictator stops acquiring information, we elicit her belief of option  $x$  being harmless, given all the information acquired. The elicitation is incentivized by the randomized Quadratic Scoring Rule adapted from Drerup *et al.* (2017) and Schlag *et al.* (2013).

In Figure 8, we present the distribution of the difference between the elicited and the rounded Bayesian belief in *Control* and *Tradeoff* respectively. The deviation of the elicited beliefs from the Bayesian beliefs does not significantly differ between *Control* and *Tradeoff* (two-sided Mann-Whitney-U test,  $p = .29$ ). In *Control*, the benchmark treatment, the dictators' elicited beliefs of  $x$  being the harmless option are on average 3.30 percentage points higher than the rounded Bayesian beliefs (two-sided Wilcoxon signed rank

Figure 8: The belief difference



This figure shows the histogram of the difference between the rounded Bayesian posterior belief and the elicited belief that option  $x$  is harmless.

test,  $p = .00$ ). In *Tradeoff*, the dictators state beliefs that are higher than the rounded Bayesian beliefs by 1.84 percentage on average (two-sided Wilcoxon signed rank test,  $p = .00$ ).

**Cognitive abilities** On average, the subjects answered 3.6 out of 5 questions in Raven’s matrices test correctly. There is no significant difference between *Control* and *Tradeoff* treatments (Chi-square  $p = .12$ ). When asked a simple question on probability, in both treatments more than 90% of the subjects answer correctly (*Treatment*: 92%, *Control*: 94%; Chi-square test  $p = .51$ ).<sup>41</sup>

<sup>41</sup>We use the following question to elicit the subjects’ understanding of probabilities: Imagine the following 4 bags with 100 fruits in each. One fruit will be randomly taken out. For which bag, the probability of taking a banana is 40%?

- A. A bag with 20 bananas.
- B. A bag with 40 bananas.
- C. A bag with 0 banana.
- D. A bag with 100 bananas.

The correct answer is B.

**SVO measure.** We elicit the social value orientation (SVO) of the subjects as a measure of their altruism. The average SVO score of all the subjects is 20.49, with no significant difference between *Tradeoff* and *Control* treatments (two-sided Mann-Whitney-U test,  $p = .84$ ). According to Murphy *et al.* (2011), 48% subjects are categorized as ‘prosocials’, 15% ‘individualists’ and 37% ‘competitive type’. The categorization is similar between *Tradeoff* and *Control*, suggesting that the treatment variation in the dictator game has no influence on the SVO measure of altruism (Chi-Square,  $p = 1.00$ ). Further, in *Tradeoff*, prosocial dictators choose the self-benefiting option  $X$  significantly less often (Chi-Square,  $p = .029$ ), while in *Control* the three categories of dictators’ decisions in the dictator game are similar (Chi-Square,  $p = .573$ ). This result suggests that the SVO measures altruistic traits that are relevant in our experimental setup.

## D Theory appendices

### D.1 Information avoidance

In this section, we discuss two additional results of our model and the respective empirical evidence: the avoidance of noisy information and of information revealing the state at once. While our experiment focuses on *noisy* information, our preference model (1) can be used to analyze both the acquisition of noisy information that arrives sequentially, as well as information that reveals the state at once. These additional theoretical results are in line with the empirical findings in this paper (see Section D.1) and also with the empirical findings by Dana *et al.* (2007) and Feiler (2014), as explained further below.

#### D.1.1 Avoidance of noisy information

Our model predicts that with or without a remunerative option, some agent types move on to the decision without acquiring any noisy information.

This may be surprising in the scenario where no option is remunerative. In this scenario, indeed, for the types with  $u' > 0$  it is optimal to acquire as much information as possible. In particular, these agents will not avoid information completely. However, other agent types have a threshold level of certainty. They are content when believing is sufficiently likely that they can spare the other from harm. Any further certainty beyond the threshold does not increase their belief-based utility. This behaviour mirrors that of satisficing as in (Simon, 1955). If, given the prior belief, they are already more certain than the threshold requires that one of the options is harmless, they are indifferent between continuing and stopping, and may as well stop.

When one of the options is remunerative, the incentives to acquire information are different. First, the agent would decide to not acquire noisy information only if she would choose the remunerative option at the prior belief. Otherwise, she would ‘fish for information’ that justifies this choice (compare to Lemma 3 and Theorem 1). Second, when considering to avoid information and to choose the remunerative option immediately or to acquire further information, the agent is aware that further information poses an undesirable risk since it might reverse her decision from the remunerative to the

non-remunerative option. When this risk outweighs her utility gain from having more certain beliefs, the agent avoids noisy information completely. The proof is in Section D.1.3.

**Proposition 2** *For any  $r = 0$  ( $r > 0$ ) and for any prior  $p_0 \in (0, 1)$ , there is a set  $S_r(p_0)$  of preference types  $u$  for which it is (strictly) optimal to avoid information completely.*

**Empirical findings.** In line with Proposition 2, in the experiment, we find that 15% of the dictators do not acquire any noisy information in the *Tradeoff–NoForce* treatment (Chi-Square  $p = 0.00$ ).<sup>42</sup> Among those, 96% choose the remunerative option  $x$  (25/26). Here, theory suggests that these dictators avoid information because they are worried about bad news arriving, indicating that  $x$  harms the other.

We find that 7% of the dictators do not acquire any noisy information in the *NoForce–Control* treatment (Chi-Square  $p = 0.00$ ). Among those, only 17% choose the remunerative option  $x$  (2/12). Here, theory suggests that the dictators who avoid noisy information completely are satisfied with 65% certainty that  $y$  is the harmless option.

### D.1.2 Avoidance of Information Revealing the State at Once

Proposition 3 shows that when there is a remunerative option, then, for any prior belief, there are types of dictators who would avoid information that reveals the state all at once. The proof is in Section D.1.4.

**Proposition 3** *1. Take  $r \geq 0$ . For any prior  $p_0 \in (0, 1)$ , there is a set of preference types  $u$  that prefer no information over receiving a signal that perfectly reveals the state.*

*2. Take any prior beliefs  $p'_0 < p_0 \in (0, 1)$ . Take a type where  $u$  is strictly concave in the second argument. If she prefers no information over receiving a signal that perfectly reveals the state when holding the prior belief  $p_0$ , she also does so when holding the larger prior belief  $p_0$ .*

---

<sup>42</sup>Recall that in Force, it is not feasible to stop immediately.

Proposition 3 is consistent with the empirical finding of Dana *et al.* (2007), who, in a dictator environment similar to ours, find that a significant fraction of dictators avoids information that reveals the *ex-ante* unknown state all at once. Feiler (2014) further documents that the fraction of dictators who avoid such perfectly revealing information increases with the dictator's prior belief that a self-benefiting option has no negative externality. The second part of Proposition 3 shows that the model predicts also this finding for a large class of preference types.

### D.1.3 Proof of Proposition 2

Take  $r = 0$ . For any prior  $p_0(0, 1)$ , there is an open set of thresholds  $l(x) \in (0, 1)$  and  $l(y) \in (0, 1)$  such that  $p_0 > l(x)$  or  $p_0 > 1 - l(y)$ . The claim for  $r = 0$  follows then from Lemma 5. Take  $r > 0$ . It follows from the characterization of the belief cutoffs after Lemma 2 that  $p_h \leq l(x)$ . Hence, if  $p_0 > l(x)$ , the agent stops acquiring information immediately in the equilibrium given by  $p_h$  and  $p_l$ .

### D.1.4 Proof of Proposition 3

The first item of Proposition 3 is a corollary of Proposition 2: it says that a preference type prefers to receive no information over all possible information acquisition strategies, including those given by the belief cutoffs  $p_l = \epsilon$  and  $p_h = 1 - \epsilon$ , which yield information arbitrarily close to a signal that reveals the state perfectly. For example, types with  $p_0 > l(x)$  cannot achieve a higher payoff than from stopping directly and choosing  $x$  since this strategy yields  $V(p_0) = r + u(x, p_0) = r$ , given the definition of  $l(x)$ , (8). The payoff when revealing the state is  $p_0V(1) + (1 - p_0)V(0) = p_0r + (1 - p_0) \max(0, r + u(x, 0))$ , which is strictly smaller than  $r$  if  $u(x, 0) < 0$ . Hence, these types strictly prefer to receive no information over receiving a signal that reveals the state perfectly.

Now, we prove the second item of Proposition 3. Fix a prior  $p_0 \in (0, 1)$ . Take any prior belief  $p'_0 < p_0 \in (0, 1)$ . Consider an agent type with  $u$  strictly concave in the second argument. Suppose that, given the prior  $p'_0$ , she prefers to receive no information over receiving a signal that perfectly reveals the state. We show that the agent also prefers to avoid information when the

prior is  $p_0$ . There are two cases. In the first case,  $x = \arg \max_{a \in \{x,y\}} U(a, 0; r)$ . This implies that the agent strictly prefers  $x$  over  $y$  at any belief  $p \in [0, 1]$ . The strict concavity of  $u$  implies

$$u(x, p_0) > p_0 u(x, 1) + (1 - p_0) u(x, 0), \quad (31)$$

which is equivalent to

$$V(p_0) > p_0 V(1) + (1 - p_0) V(0), \quad (32)$$

which shows that the agent strictly prefers to avoid information at the prior belief  $p_0$ . In the second case,  $y \in \arg \max_{a \in \{x,y\}} U(a, 0; r)$ . Since we assumed that the agent avoids information given the prior belief  $p'_0$ ,

$$V(p'_0) > (1 - p'_0) V(0) + p'_0 V(1) \quad (33)$$

Now, we use that the agent prefers  $x$  at  $p_h$  given Lemma 3, and that  $V(0) = \max_{a \in \{x,y\}} U(a, 0; r) = U(y, 0; r) = u(y, 1) = 0$ . Thus, (33) implies

$$r + u(x, p'_0) > r p'_0. \quad (34)$$

Rearranging,

$$r > -\frac{u(x, p'_0)}{1 - p'_0} \quad (35)$$

It follows from the concavity of  $u$  that  $\frac{-u(x,p)}{1-p} = \frac{u(x,1)-u(x,p)}{1-p}$  is strictly decreasing in  $p$ . Thus,

$$r > -\frac{u(x, p_0)}{1 - p_0}, \quad (36)$$

or equivalently,

$$r + u(x, p_0) < r p_0. \quad (37)$$

Thus, the type also prefers to avoid information when the prior is  $p_0$ . This finishes the proof of the second item.

## D.2 Self-image concerns

### D.2.1 Model Variation and Result

The game is as in Section 1.1: there is an agent who can observe the information process  $(Z_t)_{t \geq 0}$  at no cost. The agent can stop at any time, and chooses between  $x$  and  $y$  subsequently. Relative to Section 1.1, we modify the preferences of the agent. We follow the existing literature on self-image concerns in the prosocial domain in three aspects: first, the agent has a prosocial type  $\theta \in [0, 1]$  that captures how much she cares about the welfare of the other relative to her own remuneration. Second, she is concerned about her prosocial self-image, i.e., her utility depends on the belief about her prosocial type (e.g., Grossman and van der Weele, 2017). Third, the stopped belief together with the subsequent choice of an action  $a$  are ‘diagnostic’ about her own prosocial type (Bodner and Prelec, 2003).

Formally, a type  $\theta$  who stops observing the process at  $t = \tau$ , holding a belief  $p = p_\tau$ , and who chooses  $a \in \{x, y\}$ , has the following utility (compare to (1)),

$$U(a, p; \theta, r) = \begin{cases} u(x, p) + rw(\theta) & \text{if } a = x, \\ u(y, 1 - p) & \text{if } a = y, \end{cases} \quad (38)$$

where the belief-based utility is

$$u(a, q) = \begin{cases} q + \psi E(\theta | p_\tau = q, a = x) & \text{if } a = x, \\ q + \psi E(\theta | p_\tau = 1 - q, a = y) & \text{if } a = y. \end{cases} \quad (39)$$

The first term in (39) represents the prosocial concern and the second term represents the diagnostic part of the utility, with  $\psi > 0$ . The function  $w$  is continuously differentiable and strictly decreasing in the prosocial type, so that agents with higher prosocial type care less about their own remuneration from choosing the self-benefiting option  $x$ . A strategy is a mapping  $\sigma$  from types to stopping times.

We study the situation with egoistic motive ( $r = 1$ ). The following result provides sufficient conditions for the existence of a perfect Bayesian equilibrium that is *monotone*, meaning that more prosocial types acquire more

information about the state in a Blackwell sense.<sup>43</sup>

**Theorem 2** *Let  $r = 1$ . Take a distribution  $F$  of the prosocial types  $\theta$  with log-concave density and support  $[0, 1]$ . Let  $0 < \psi(\psi + w') < 1$ ,  $w(1) < 1 + \psi < w(0)$ , and  $p_0 + \psi E(\theta) + (1 - p_0)w(1) > (1 + \psi)$ . For  $\epsilon > 0$  sufficiently small, there exists a monotone perfect Bayesian equilibrium with the following properties.*

1. *There is  $k \in (0, 1)$  such that all types  $\theta < k$  acquire no information and all types  $\theta \geq k$  acquire partial information about the state.*
2. *If a type  $\theta \geq k$  is more prosocial than another type  $\theta'$ , i.e.,  $\theta > \theta'$ , then, the type  $\theta$  uses a strategy that generates signals that are Blackwell strictly more informative about the state relative to the strategy of  $\theta'$ .*

The conditions in the statement of Theorem 2 are sufficient, but not necessary. The following proof shows how the conditions ensure that the equilibrium is interior (i.e.,  $k \in (0, 1)$ ), monotone, and that the best response of the types  $\theta \geq k$  is pinned down locally by a first-order condition.

## D.2.2 Proof: candidate equilibria

**Strategies.** Let  $\epsilon > 0$ . Take  $k \in [0, 1]$ . We define an associated strategy  $\sigma^k$ . Take a type  $\theta < k$ . The type stops observing the information process directly at the prior  $p_0$ . Take a type  $\theta \geq k$ . The type uses the belief cutoff  $p_l = \epsilon$ , that is, she stops observing the information process when  $p_t \leq p_l$ , and chooses  $y$  subsequently. She also stops if  $p_t \geq v^k(\theta)$  and chooses  $x$  subsequently, where  $v^k$  is defined as follows,

$$v^k(\theta) = c_1 e^{\int_k^\theta g^k(z) dz} + \epsilon \quad (40)$$

for  $g^k(z) = \psi \left[ \psi z - \psi E(\theta | \theta \geq k) + w(z) - 1 + 2\epsilon \right]^{-1}$  and for a constant  $c_1 \geq p_0$ , so that  $v^k(k) \geq p_0 + \epsilon$  for all  $k$ . We see that  $v^k(\theta)$  solves the following ordinary differential equation (ODE).

$$(v^k)'(\theta) = v^k(\theta)g(\theta). \quad (41)$$

---

<sup>43</sup>There are also pooling equilibria. However, the goal here is to show that there may exist equilibria in which types self-signal about their prosociality by stopping at more or less informative beliefs about the consequences of their action choice on others.

Later, when analyzing the best response, we will show that the ordinary differential equation (41) will turn up as first-order condition of the agent's maximization problem.

**Feasibility.** In the rest of this section, we restrict the domain of  $k$  to ensure that  $v^k(\theta) \leq 1 - \epsilon$  for all  $\theta \geq k$ , so that  $\sigma^k$  is a feasible strategy.

Recall that  $w$  is continuous, and so is  $E(\theta|\theta \geq k)$  in  $k \in [0, 1]$  since the distribution of  $\theta$  has a density and since its support is  $[0, 1]$ . The assumptions  $w(1) < 1 + \psi < w(0)$  imply

$$w(1) < 1 + \psi E(\theta|\theta \geq 1) \quad (42)$$

and

$$w(0) > 1 + \psi E(\theta|\theta \geq 0). \quad (43)$$

Altogether, the intermediate value theorem implies that there is  $\tilde{k} \in (0, 1)$  such that

$$w(\tilde{k}) = 1 + \psi E(\theta|\theta \geq \tilde{k}). \quad (44)$$

Recall that  $w$  is strictly decreasing in  $\theta$  and note that  $E(\theta|\theta \geq k)$  is strictly increasing in  $k$  since the density of  $\theta$ 's distribution is log-concave (see e.g., Burdett, 1996). Hence,

$$w(\theta) < 1 + \psi E(\theta|\theta \geq k) \quad (45)$$

for all  $k \geq \tilde{k}$  and all  $\theta \geq k$ .

The following Lemma 6 shows that  $g^k > 0$ . Further, the lemma shows that  $g$  is uniformly bounded above. Using, (41), this implies that for any  $c_1 \geq p_0$ , there is minimal  $k(c_1) \in [\tilde{k}, 1]$  such that  $p_0 + \epsilon \leq v^k(\theta) \leq 1 - \epsilon$  for all  $k \geq k(c_1)$  and  $\theta \geq k$ . Hence,  $\sigma^k$  is a well-defined strategy for  $k \geq k(c_1)$ . In the remainder of this proof, let  $k \geq k(c_1)$ .

**Lemma 6** *Let  $w(0) > 1 + \psi$ , and  $0 < \psi(\psi + w') < 1$ . Then,  $g^k(\theta) > 0$ ,  $(v^k)'(\theta) > 0$  and  $(v^k)''(\theta) > 0$  for all  $k, \theta \in [0, 1]$ . Further,  $g^k(\theta)$  is uniformly bounded.*

**Proof.** Recall (41). Hence,  $(v^k)' > 0$  if  $g > 0$ . Note that  $g^k(\theta) > 0$  if

$$\psi\theta + w(\theta) \geq 1 + \psi\mathbf{E}(\theta|\theta \geq k) \quad (46)$$

since  $\epsilon > 0$ . Now, (46) holds for  $\theta = 0$  since  $w(0) \geq 1 + \psi$ . Further, it holds for  $\theta > 0$  since  $\psi + w' > 0$  implies  $\psi\theta + w(\theta) > 1 + \psi$ . Second, note that it follows from (41) that

$$(v^k)'' = (v^k)'g + vg' = (v^k)(g^2 + g') \quad (47)$$

Hence,  $(v^k)'' > 0$  if  $g^2 + g' > 0$ . Let  $h = \frac{1}{g}$ , i.e.,  $g = \frac{1}{h}$ . Then,  $g' = -\frac{h'}{h^2}$ . Hence,  $g^2 + g' > 0$  if  $1 - h' > 0$ . By definition,  $h' = \psi(\psi + w')$ , and, by assumption,  $\psi(\psi + w') < 1$ . We conclude that  $(v^k)'' > 0$ . Third, note that the assumption that  $h' = \psi(\psi + w') > 0$  implies that  $g$  is maximal at  $\theta = 0$ . This implies  $g^k(\theta) \leq \frac{1}{\epsilon + w(0) - 1 - \psi}$  for all  $k, \theta \in [0, 1]$ . Since also  $g > 0$  for all  $k, \theta \in [0, 1]$ , we conclude that  $g$  is uniformly bounded. ■

**Beliefs.** We define a belief system  $\mu^k(a, q)$  for  $k \geq k(c_1)$ . Given  $\sigma^k$ , the following beliefs are well-defined by Bayes-consistency:

$$\mathbf{E}(\theta|p_\tau = q, a; \sigma^k) = \begin{cases} (v^k)^{-1}(q) & \text{for } q \geq v^k(k) \text{ and } a = x, \\ \mathbf{E}(\theta|\theta \leq k) & \text{for } q = p_0 \text{ and } a = x, \\ \mathbf{E}(\theta|\theta \geq k) & \text{for } q = \epsilon \text{ and } a = y, \end{cases} \quad (48)$$

where we used that  $(v^k)' > 0$  from Lemma 6. Let

$$\mu(y, q) = \begin{cases} \frac{q}{1-\epsilon}\mathbf{E}(\theta|\theta \geq k) & \text{for } q \leq 1 - \epsilon, \\ \mathbf{E}(\theta|\theta \geq k) & \text{for } q \geq 1 - \epsilon, \end{cases} \quad (49)$$

and

$$\mu(x, q) = \begin{cases} 0 & \text{for } q \leq p_0 - \epsilon, \\ \frac{q - (p_0 - \epsilon)}{\epsilon}\mathbf{E}(\theta|\theta \leq k) & \text{for } p_0 - \epsilon \leq q < p_0, \\ \mathbf{E}(\theta|\theta \leq k) & \text{for } p_0 \leq q < v^k(k), \\ (v^k)^{-1}(q) & \text{for } v^k(k) \leq q \leq v^k(1), \\ 1 & \text{for } q > v^k(1). \end{cases} \quad (50)$$

We see that  $\mu(a, q)$  is consistent with the beliefs that are implied by  $\sigma^k$  through Bayes rule, (48). We see that  $\mu(x, q)$  is weakly increasing in  $q$  and continuously differentiable at all  $q > v^k(k)$ . The pairs  $(\sigma^k, \mu^k)$  are the candidate equilibria.

### D.2.3 Proof: best response to a candidate equilibrium

Take  $\sigma^k$  and  $\mu^k$  and consider the best response. Consider the candidate belief-based utility

$$u(a, q) = q + \psi\mu(a, q). \quad (51)$$

Take any  $\theta \geq k$  for which it is not optimal to stop directly at the prior given  $\sigma^k$  and  $\mu^k$ . In the following, we rely on the previous analysis of the model with belief-based utility from Section 1.1 to characterize the best response for  $\theta$ .

**Claim 4** *There is a best response characterized by belief cutoffs  $p_l < p_0$  and  $p_h > p_0$ , meaning that the agent stops if and only if  $p_t \leq p_l$  or  $p_t \geq p_h$ .*

Let  $V(p) = \max_{a \in \{x, y\}} U(a, p; \theta, r)$ . Let  $\bar{V}$  be the smallest concave function with  $\bar{V}(p) \geq V(p)$  for all  $p \in [\epsilon, 1 - \epsilon]$ . Following the steps of the proofs of Lemma 1 and Lemma 2 in Appendix B verbatim, we show that the results of the two lemmas hold. In particular, the argument from the proof of Lemma 2 shows the claim and that

$$V(p_l) = \bar{V}(p_l), \quad (52)$$

$$V(p_h) = \bar{V}(p_h), \quad (53)$$

compare to (21).

**Claim 5** *It is optimal to choose  $x$  at  $p_h$ .*

Compare the statement of the claim to Lemma 3. We note that  $u(y, 1) \leq u(x, 1)$  given  $\sigma^k$  and  $\mu^k$ , and normalize the belief-based utility to  $\tilde{u}(a, q) = u(a, q) - u(x, 1)$ . This way,  $\tilde{u}(x, 1) = 0$  and  $\tilde{u}(y, 1) \leq 0$ . Following the steps of the proof of Lemma 3 verbatim, while replacing  $u$  by  $\tilde{u}$  therein, shows that

it is optimal to choose  $x$  at  $p_h$  and that

$$V(p_h) > V(p_l). \quad (54)$$

**Claim 6** *It is optimal to choose  $y$  at  $p_l$  and it holds that  $p_l = \epsilon$  when  $\epsilon > 0$  is sufficiently small.*

Compare the statement of the claim to Lemma 4. First, we argue that it is *not* optimal to choose  $x$  at  $p_l$ . This implies that it is also optimal to choose  $x$  at any  $p \geq p_l$  given the definitions of the payoffs (38) and since the belief utility  $q + \mu(x, q)$  is increasing in  $q$  and the belief utility  $\mu(y, q)$  is decreasing in  $q$ . Hence,  $V(p) = \max_{a \in \{x, y\}} U(a, p; \theta, r) = U(x, p; \theta, r)$  for  $p \geq p_l$ . Given the definition of  $\mu(x, q)$ , (50), we note that  $\bar{V}(p_0) = V(p_0)$ , so that it is optimal for  $\theta$  to stop acquiring information directly. This contradicts with the initial assumption. Now, following the lines of the proof of Lemma 4 and using (54) and that  $u(y, q)$  is weakly increasing in  $q$  for  $q > 1 - p_0$ , we show that  $p_l = \epsilon$ , that is Lemma 4 holds.

Together with Claim 4 - 6, the next result, Claim 7 finishes the characterization of the best response when it is not optimal to stop the information process directly.

**Claim 7** *We have  $p_h = v^k(\theta)$ .*

**Proof.** Take any  $\theta \geq k$  for which it is not optimal to stop directly at the prior. Given Claim 4 - 6, it is optimal for  $\theta$  to use a strategy with belief cutoffs  $p_l = \epsilon$  and  $q$  for some  $q > p_0$ : she stops only if  $p_t \leq \epsilon$ , after which she chooses  $y$ , or she stops if  $p_t \geq q$ , in which case she chooses  $x$ .

Now, we show that the optimal upper belief cutoff is given by  $q = v^k(\theta)$ . For this, fix  $p_l = \epsilon$ , and consider the agent's expected payoff from using the upper cutoff  $q \geq p_0$ ,

$$\frac{q - p_0}{q - \epsilon} u(y, 1 - \epsilon) + \frac{p_0 - \epsilon}{q - \epsilon} \left[ u(x, q) + w(\theta) \right]. \quad (55)$$

It turns out to be useful for the algebra to subtract  $u(y, 1 - \epsilon)$  from the objective function (55). Doing so, and taking the first-order condition with

respect to  $q$  gives

$$\begin{aligned} & \frac{p_0 - \epsilon}{q - \epsilon} u'(x, q) - \frac{p_0 - \epsilon}{(q - \epsilon)^2} (u(x, q) + w(\theta) - u(y, 1 - \epsilon)) = 0 \\ \Leftrightarrow & u'(x, q)(q - \epsilon) = u(x, q) + w(\theta) - u(y, 1 - \epsilon). \end{aligned} \quad (56)$$

Recalling the definition of  $\mu^k$ , (50),

$$u(x, q) = q + \psi(v^k)^{-1}(q) \quad (57)$$

for  $q \geq v^k(k)$ . Combining (56) and (57),

$$q - \epsilon + \psi((v^k)^{-1})'(q)(q - \epsilon) = q + \psi(v^k)^{-1}(q) + w(\theta) - u(y, 1 - \epsilon). \quad (58)$$

Now, we would like to check if  $q = v^k(\theta)$  satisfies the first-order condition.

For  $q = v^k(\theta)$ , we have  $((v^k)^{-1})'(q) = \frac{1}{(v^k)'(\theta)}$ . Rewriting (58),

$$-\epsilon + \psi \frac{v^k(\theta) - \epsilon}{(v^k)'(\theta)} = \psi\theta + w(\theta) - u(y, 1 - \epsilon).$$

Using that  $u(y, 1 - \epsilon) = 1 - \epsilon + \psi E(\theta | \theta \geq k)$ ,

$$\begin{aligned} & -\epsilon + \psi \frac{v^k(\theta) - \epsilon}{(v^k)'(\theta)} = \psi\theta + w(\theta) - 1 + \epsilon - \psi E(\theta | \theta \geq k) \\ \Leftrightarrow & (v^k(\theta) - \epsilon)g(\theta) = (v^k)'(\theta) \end{aligned} \quad (59)$$

for  $g(z) = \psi \left[ \psi z - \psi E(\theta | \theta \geq k) + w(z) - 1 + 2\epsilon \right]^{-1}$ . Using the definition (40), one checks that  $v^k(\theta)$  satisfies the ordinary differential equation (59), hence the first-order condition (56).

To conclude the proof, we show that the solution to the first-order condition (56) is unique and maximizes (55). For this, we argue that  $u''(x, q) < 0$  for  $q \geq v^k(k)$ .<sup>44</sup> Recall that  $u(x, q) = q + (v^k)^{-1}(q)$  for  $q \geq v^k(k)$ . Recall from Lemma 6 that  $(v^k)''(\theta) > 0$ , so  $v^k$  is strictly convex. Since  $v^k$  is strictly increasing in  $\theta$  by Lemma 6, this implies that  $(v^k)^{-1}(q)$  is strictly concave, hence  $u(x, q)$  is strictly concave for  $q \geq v^k(k)$ . The concavity of  $u(x, q)$  implies

---

<sup>44</sup>For  $q = v^k(k)$ ,  $u''(x, q)$  is the right-sided derivative.

that

$$u'(x, q)(q - \epsilon) - \left[ u(x, q) - w(\theta) + u(y, 1) \right] \quad (60)$$

is strictly decreasing in  $q$  for  $q \geq v^k(k)$  when  $\epsilon \approx 0$ .<sup>45</sup> Thus, using (56), the objective function (55) is strictly concave for  $q \geq v^k(k)$ . Therefore,  $v^k(\theta)$  is the unique solution to (56) and maximizes (55) across all  $q \geq v^k(k)$ .

Finally, note that we constructed  $\mu(x, q)$  and thereby  $u(x, q)$  such that (55) is not maximized by any  $p_0 \leq q < v^k(k)$ . Suppose otherwise. Then, the strategy given by the belief cutoffs  $p_l = \epsilon$  and  $q$  must maximize  $E(V(p_\tau))$  with  $V$  defined as after Claim 4. Note that, given the definitions of  $\mu(x, q)$ , (50), and  $U(a, p; \theta, r)$ , (38), the function  $V$  is linear on  $[p_0, v^k(k))$  and has a jump upwards at  $v^k(k)$ . This shows that, fixing  $p_l = \epsilon$ , the expected payoff  $E(V(p_\tau))$  is maximized by either  $q = p_0$  or by  $q = v^k(k)$  on the interval  $[p_0, v^k(k)]$ . Our initial assumption that it is not optimal to stop acquiring information directly rules out  $q = p_0$ . We conclude that  $v^k(\theta)$  is the unique global maximizer. ■

#### D.2.4 Proof: equilibrium construction

We define an auxiliary map  $f$  from  $c_1$  to a type  $f(c_1)$ . This map will split the types into those that acquire no information under the best response to  $(\sigma^k, \mu^k)$ , and those that acquire some information. First, we restrict the domain. Note that for  $c_1 \geq p_0$ , we have

$$v^k(k) \geq p_0 + \epsilon, \quad (61)$$

given (40). Given Lemma 6, this implies  $v^k(1) > p_0$ . Take  $\delta > 0$ . Denote  $\bar{c}_1 = 1 - \delta - \epsilon$ , which implies that, given  $c_1 = \bar{c}_1$ , we have  $v^k(k) = 1 - \delta$ . For any  $c_1$ , take  $k = k(c_1)$ , meaning that, given  $c_1$ , it holds

$$v^{k(c_1)}(1) = 1 - \epsilon. \quad (62)$$

The inequalities (61) and (62) together with  $(v^k)' > 0$  (see Lemma 6) imply that for all  $p_0 \leq c_1 \leq \bar{c}_1$  and  $\theta \geq k(c_1)$ , it holds  $v^{k(c_1)}(\theta) \in [p_0 + \epsilon, 1 - \epsilon]$ . Hence, the strategy  $\sigma^{k(c_1)}$  is well-defined when  $p_0 \leq c_1 \leq \bar{c}_1$ .

---

<sup>45</sup>To see why, consider the derivative with respect to  $q$  for  $\epsilon \approx 0$ . It is given by  $u''(x, q)q + u'(x, q) - u'(x, q) = u''(x, q) < 0$ .

We consider a candidate best response to  $(\sigma^{k(c_1)}, \mu^{k(c_1)})$  where the agent stops only if  $p_t \leq \epsilon$ , or if  $p_t \geq q$ ; further, if she stops at  $p_t \leq \epsilon$ , she chooses  $y$ , and if she stops above  $q$ , then she chooses  $x$  (compare to the results in Section D.2.3.) Given Claim 7, it is sufficient to restrict to candidate strategies with  $q \geq p_0 + \epsilon$  since  $v^k(\theta) \geq p_0 + \epsilon$ , see (61). Another candidate best response is that the agent stops observing the information process directly. Given Lemma 3, she chooses  $x$  when stopping directly. Given  $\sigma^{k(c_1)}$  and  $\mu^{k(c_1)}$ , a type  $\theta$  prefers to stop observing the information process directly to using the candidate strategy given by the belief cutoffs  $\epsilon$  and  $q \geq p_0 + \epsilon$  if

$$\begin{aligned} & \frac{q - p_0}{q - \epsilon} \left[ (1 - \epsilon) + \psi \mu(y, 1 - \epsilon) \right] + \frac{p_0 - \epsilon}{q - \epsilon} \left[ q + \psi \mu(x, q) + w(\theta) \right] \\ & \leq p_0 + \psi \mathbb{E}(\theta | \theta \leq k(c_1)) + w(\theta). \end{aligned} \quad (63)$$

Since  $w'(\theta) < 0$  and since  $q > p_0$ , the difference of the left hand side minus the right hand side is strictly increasing in  $\theta$ . Hence, there exists a unique cutoff  $\theta(q) \in [0, 1]$  such that a type  $\theta$  prefers to stop directly if  $\theta < \theta(q)$  and a type prefers to use  $p_l = \epsilon$  and  $q$  if  $\theta > \theta(q)$ . Since  $q$  was arbitrary, we obtain a cutoff

$$f(c_1) = \inf \{ \theta(q) : q \geq p_0 + \epsilon \} \in [0, 1] \quad (64)$$

such that a type  $\theta$  prefers to stop directly over *all* strategies given by belief cutoffs  $\epsilon$  and some  $q \geq p_0 + \epsilon$  if  $\theta < f(c_1)$ , and prefers some such strategy given by belief cutoffs  $\epsilon$  and  $q$  over stopping directly if  $\theta > f(c_1)$ .

**Construction.** We claim that there is  $c_1 \in [p_0, \bar{c}_1]$  such that

$$f(c_1) = k(c_1). \quad (65)$$

For such  $c_1$ , the characterization of the best response in Section D.2.3 shows that, for all types  $\theta < k(c_1)$  it is strictly optimal to acquire no information, and for all types  $\theta \geq k(c_1)$  it is optimal to acquire some information and to use the strategy with the belief cutoffs  $p_l = \epsilon$  and  $p_h = v^{k(c_1)}(\theta)$ . Since  $\mu^{k(c_1)}$  is a belief system that is consistent with  $\sigma^{k(c_1)}$  by definition, this shows that  $\mu^{k(c_1)}$  and  $\sigma^{k(c_1)}$  constitute a perfect Bayesian equilibrium, and finishes the proof of Theorem 1.

To show (65), it is useful to define by  $\phi(c_1)$  the utility difference of the type  $\theta = k(c_1)$  from the strategy where she stops directly at the prior and chooses  $x$  relative to the strategy where she stops only if  $p_t \leq \epsilon$ , choosing  $y$  thereafter, or if  $p_t \geq v^{k(c_1)}(k(c_1))$ , choosing  $x$  thereafter.

**Step 1** *If  $c_1 = p_0$  and  $\epsilon > 0$  is sufficiently small, it is not optimal for the type  $\theta = k(c_1)$  to stop directly given  $\sigma^{k(c_1)}$  and  $\mu^{k(c_1)}$ . That is,  $\phi(p_0) < 0$ .*

**Proof.** Note that by definition of  $\mu^k$  ((48)),

$$\mu^{k(c_1)}(x, q) = \begin{cases} \mathbb{E}(\theta | \theta \leq k(c_1)) & \text{for } q = p_0, \\ (v^{k(c_1)})^{-1}(q) & \text{for } q \geq v^{k(c_1)}(k(c_1)). \end{cases} \quad (66)$$

Recall that, for  $c_1 = p_0$ , we have  $v^{k(c_1)}(k(c_1)) = p_0 + \epsilon$ , given (40). Therefore,  $(v^{k(c_1)})^{-1}(q)$  equals  $k(c_1)$  for  $q = p_0 + \epsilon$ , which is strictly larger than  $\mathbb{E}(\theta | \theta \leq k(c_1))$ . When  $\epsilon \approx 0$ ,  $\theta$  is therefore strictly better off when observing the information process as long as  $p_t \in [p_0, p_0 + \epsilon]$ , choosing  $x$  when stopping at  $p_t = p_0 + \epsilon$ , and choosing  $y$  when stopping at  $p_t = p_0$  than when using the strategy where she stops directly at  $p_0$  and choosing  $x$ . Hence, it is not optimal to stop at the prior directly. ■

**Step 2** *If  $c_1 = \bar{c}_1$ , and  $\epsilon, \delta > 0$  are sufficiently small, it is strictly optimal for the type  $\theta = k(c_1)$  to stop directly given  $\sigma^{k(c_1)}$  and  $\mu^{k(c_1)}$ . That is,  $\phi(\bar{c}_1) > 0$ .*

**Proof.** Recall that  $\bar{c}_1 = 1 - \delta - \epsilon$ . It follows from Claim 5 - Claim 7 that it is either optimal for  $\theta = k(c_1)$  to stop directly or to use the strategy given by the belief cutoffs  $p_l = \epsilon$  and  $p_h = v^{k(c_1)}(k(c_1))$ . Let  $k = k(c_1)$ . Stopping directly is strictly optimal if

$$\begin{aligned} U(x, p_0; \theta, r) &> \frac{p_0 - \epsilon}{p_h - \epsilon} U(x, p_h; \theta, r) + \frac{p_h - p_0}{p_h - \epsilon} U(y, \epsilon; \theta, r), \\ \Leftrightarrow p_0 + \psi \mathbb{E}(\theta | \theta \leq k) + w(\theta) & \quad (67) \\ &> \frac{p_0 - \epsilon}{v^k(k) - \epsilon} \left[ v^k(k) + \psi k + w(\theta) \right] + \frac{v^k(k) - p_0}{v^k(k) - \epsilon} \left[ 1 - \epsilon + \psi \mathbb{E}(\theta | \theta \geq k) \right]. \end{aligned}$$

Note that  $\delta \approx 0$  and  $\epsilon \approx 0$  implies  $\bar{c}_1 = 1 - \delta - \epsilon \approx 1$ , which in turn implies  $k(c_1) \approx 1$  by the definition of  $k(c_1)$ . The equation (67) holds for  $\delta > 0$ ,  $\epsilon > 0$ ,

and  $k(c_1) - 1$  sufficiently small if

$$\begin{aligned} p_0 + \psi \mathbb{E}(\theta | \theta \leq k(c_1)) + w(1) &> p_0(1 + \psi + w(1)) + (1 - p_0)(1 + \psi) \\ \Leftrightarrow p_0 + \psi \mathbb{E}(\theta | \theta \leq k(c_1)) + (1 - p_0)w(1) &> 1 + \psi. \end{aligned} \quad (68)$$

The condition (68) holds because of the assumption  $p_0 + \psi \mathbb{E}(\theta) + (1 - p_0)w(1) > 1 + \psi$  of Theorem 2. Finally, since all previous arguments were true for all  $\delta > 0$  and  $\epsilon > 0$  sufficiently small, we can choose these parameters small enough so that (67) holds. This finishes the proof that  $\phi(\bar{c}_1) > 0$ . ■

Finally, note that  $\phi$  is continuous in  $c_1$ : this is because it depends on  $c_1$  only through  $k(c_1)$  and  $v^{k(c_1)}(k(c_1))$ . However,  $k(c_1)$  is continuous in  $c_1$  and  $v^{k(c_1)}(k(c_1)) = c_1 + \epsilon$  is also continuous in  $c_1$ . Using Step 1 and Step 2, it follows from the intermediate value theorem, that there is  $c_1 \in (p_0, \bar{c}_1)$  such that  $\phi(c_1) = 0$ . This implies  $f(c_1) = k(c_1)$  by the definition of  $f$  and  $\phi$ .

# Online appendix: experimental instructions

## E Original instructions in German

In this online appendix, we include the original instructions that we used in the experiment. The original instructions are paper-based and in German language. They are similar in all four treatment: *NoForce Tradeoff*, *Force-Tradeoff*, *Force-Control* and *NoForce-Control*. We include the instructions in *NoForce-Tradeoff* in full and point out the deviation from them in the three other treatments respectively. We include the English translation of these instructions in Online Appendix F.

### E.1 Treatment: NoForce-Tradeoff

#### Allgemeine Erklärungen

Wir begrüßen Sie zu dieser Studie! Im Rahmen dieser Studie können Sie eine nicht unerhebliche Summe Geld verdienen. Lesen Sie die folgenden Erklärungen daher bitte gründlich durch! Wenn Sie Fragen haben, strecken Sie bitte Ihre Hand aus der Kabine – wir kommen dann zu Ihrem Platz.

**Während der Studie ist es nicht erlaubt, mit den anderen Studienteilnehmern zu sprechen, Mobiltelefone zu benutzen oder andere Programme auf dem Computer zu starten.** Die Nichtbeachtung dieser Regeln führt zum Ausschluss aus der Studie und von allen Zahlungen. Ihr Einkommen aus dieser Studie bekommen Sie am Ende der Studie bar ausbezahlt. Während der Studie sprechen wir nicht von Euro, sondern von Punkten. Ihre gesamte Auszahlung wird also zunächst in Punkten berechnet und dann am Ende in Euro umgerechnet, wobei gilt:

**1 Punkt = 5 Cent.**

**Teilnehmerzuordnung:** Durch eine Zufallsentscheidung hat Ihnen der Com-

puter aus allen Studienteilnehmern in diesem Raum einen anderen Teilnehmer zugeordnet. Im folgenden bezeichnen wir den Ihnen zugeordneten Studienteilnehmer als 'den anderen Teilnehmer'.

Benutzen Sie diesen Teilnehmerbogen gerne als Referenz während der Bearbeitung am Computer. Im Vorlauf zur Bearbeitung der Studie werden wir Sie bitten, einige Kontrollfragen zu bearbeiten.

## Wahrscheinlichkeiten

In diesem Abschnitt möchten wir Sie ein wenig mit mathematischen Wahrscheinlichkeiten vertraut machen.

Stellen Sie sich folgende Situation vor: In einem Raum sitzen 10 Teilnehmer. Jeder Teilnehmer hat eine Box. Die Teilnehmer können keine der Boxen sehen, wissen aber folgendes:

- In jeder Box befinden sich 10 Bälle.
- 5 der 10 Teilnehmer haben eine Box mit 6 weißen Bällen und 4 schwarzen Bälle (Situation A).
- 5 der 10 Teilnehmer haben eine Box mit 4 weiße Bällen und 6 schwarze Bällen (Situation B).

Stellen Sie sich nun folgendes vor: Sie sind einer der Teilnehmer. Die Wahrscheinlichkeit, dass Sie sich in Situation A befinden, ist also 50 %. Ein Computer zieht einen Ball aus Ihrer Box und legt ihn nach dem Ziehen wieder zurück. Wenn Sie die Farbe des gezogenen Balles erfahren, hilft dies, besser einzuschätzen in welcher Situation Sie sich befinden.

Frage: Was ist die Wahrscheinlichkeit, dass Sie sich in Situation A befinden, gegeben dass der gezogene Ball weiß ist?

Wahrscheinlichkeiten, welche zusätzliche Information - wie die Farbe des gezogenen Balles - berücksichtigen, werden auf statistisch korrekte Weise nach

einem mathematischen Gesetz, dem Satz von Bayes' berechnet. Die Berechnung ist kompliziert, und benötigt mehrere Rechenschritte. **Daher zeigen wir Ihnen die statistisch korrekten Wahrscheinlichkeiten an, wann immer sie entscheidungsrelevant sind.**

Beispiel (Antwort zur Frage):

$$\begin{aligned} & \text{Wahrscheinlichkeit von Situation A, gegeben, dass der gezogene Ball weiß ist} \\ = & \text{W'keit von Situation A} \cdot \left[ \frac{\text{W'keit, dass der gezogene Ball in Situation A weiß ist}}{\text{W'keit, dass der gezogene Ball weiß ist}} \right] \\ = & 50\% \cdot \frac{60\%}{50\%} \\ = & 60\%. \end{aligned}$$

Die Wahrscheinlichkeit von Situation A, gegeben, dass der gezogene Ball schwarz ist, wird auf ähnliche Weise berechnet. **Bitten bearbeiten Sie nun die erste Kontrollfrage am Computer.**

## Ihre Entscheidungen

Sie erhalten nun 100 Punkte auf ihr Punktekonto. Als Nächstes fällen Sie eine Entscheidung, die Ihre eigene Auszahlungshöhe und die Auszahlungshöhe des anderen Teilnehmers beeinflusst, welcher ebenfalls 100 Punkte auf sein Punktekonto erhalten hat. Die Entscheidungen des anderen Teilnehmers haben jedoch keine Auswirkung auf Ihre Auszahlungshöhe.

Ihre Entscheidung besteht daraus, zwischen zwei Optionen, X und Y, zu wählen.

1. Eine dieser Optionen ist für den anderen Studienteilnehmer 'schädlich' und führt dazu, dass er 80 Punkte weniger ausgezahlt bekommt.
2. Die andere der Optionen hat keine Auswirkung auf den anderen Studienteilnehmer, diese Option ist 'sicher'.
3. Für je 7 von 20 Entscheidern ist die Option X sicher und Option Y schädlich (35 % Wahrscheinlichkeit), und für je 13 von 20 Entscheidern

ist die Option Y sicher und Option X schädlich (65 % Wahrscheinlichkeit). Es wurde vom Computer bereits zufällig ausgewählt, welche Option in Ihrem Fall sicher und welche schädlich ist.

4. Unabhängig davon, ob Option X oder Option Y für den anderen Teilnehmer sicher ist, bekommen Sie selber 25 Punkte mehr ausgezahlt, wenn Sie sich für Option X entscheiden.

## Ihre Information

Bevor Sie sich entscheiden, können Sie zusätzliche Information darüber erhalten, welche der Optionen der Computer für Sie als sicher ausgewählt hat. Wenn Sie die Entscheidung ohne zusätzliche Information treffen möchten, drücken Sie bitte direkt auf ‘Entscheidung ohne zusätzliche Information’. Wenn Sie die Entscheidung mit zusätzlicher Information treffen möchten, drücken Sie bitte auf ‘Entscheidung mit zusätzlicher Information’

In einer Box befinden sich 100 weiße oder schwarze Bälle. Falls Option X sicher ist, befinden sich in der Box 60 weiße Bälle, und 40 schwarze Bälle. Falls Option Y sicher ist, befinden sich in der Box 40 weiße Bälle, und 60 schwarze Bälle.

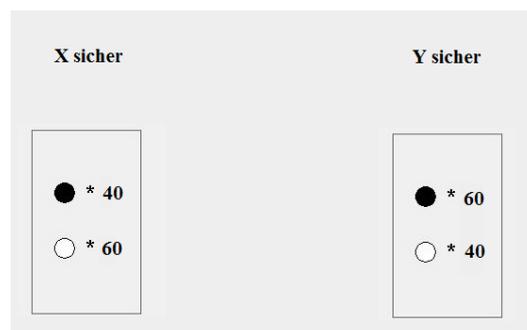


Figure 9: Box mit 40 schwarzen und 60 weißen Bällen (X sicher); Box mit 60 schwarzen und 40 weißen Bällen (Y sicher)

Sie werden die Box nicht gezeigt bekommen, aber sie können den Computer einen Ball zufällig aus der Box ziehen lassen. Dafür klicken Sie auf die Schaltfläche ‘Ein weiterer Ball’. Nach Ihrem Klick wird der Ball, den der

Computer aus der Box gezogen hat, auf dem Bildschirm eingeblendet. Danach legt der Computer den Ball in die Box zurück, sodass sich wieder 100 Bälle in der Box befinden.



Figure 10: Beispiel: Der vom Computer gezogene Ball ist weiß.

Sie können danach einen weiteren Ball ziehen lassen. Dafür klicken Sie wiederum auf die Schaltfläche 'Ein weiterer Ball'. Sie können unbegrenzt viele Bälle ziehen lassen. Wenn Sie keine Bälle mehr ziehen lassen möchten, klicken Sie auf die Schaltfläche 'Zur Entscheidung'.

Erinnern Sie sich, dass mehr weiße Bälle in der Box sind, wenn Option X sicher ist (60 weiße Bälle), als wenn Option Y sicher ist (40 weiße Bälle)? Wenn Sie also sehen, dass ein weißer Ball gezogen wurde, ist dies ein Hinweis darauf, dass Option X sicher ist. Nach jedem Ball wird die statistisch korrekte Wahrscheinlichkeit, dass Option X sicher ist, eingeblendet. Dabei werden alle Bälle, die Sie bereits gezogen haben, berücksichtigt.<sup>46</sup>

Die betroffene Person weiß nicht, ob oder wieviel Sie sich informiert haben. Sie erfahren am Ende der Studie nicht, ob X oder Y die schädliche Aktion ist. Es ist kostenlos für Sie, zusätzliche Information vor Ihrer Entscheidung zwischen X und Y zu erhalten.

<sup>46</sup>Für die genaue Berechnung der Wahrscheinlichkeit, dass Option X sicher ist, wird ein mathematisches Gesetz, der Satz von Bayes, benutzt. Die exakte Berechnungsformel, wenn z.B. ein einziger weißer Ball gezogen wurde, ist folgende:

$$\begin{aligned}
 & \text{Wahrscheinlichkeit, dass Option X sicher ist, gegeben, dass der gezogene Ball weiß ist} \\
 = & \frac{\text{W'keit, dass Option X sicher ist} \cdot \text{W'keit, dass der gezogene Ball weiß ist, wenn Option X sicher ist}}{\text{W'keit, dass der gezogene Ball weiß ist}}.
 \end{aligned}$$

Dies ist die einzige objektiv richtige Berechnungsweise.

## **E.2 Treatment: NoForce-Control**

...

### **Ihre Entscheidungen**

...

Ihre Entscheidung besteht daraus, zwischen zwei Optionen, X und Y, zu wählen.

...

4. Ihre eigene Auszahlungshöhe ist unabhängig davon, ob Sie sich für X oder Y entscheiden. Ihre Auszahlungshöhe ist auch unabhängig davon, ob Option X oder Option Y für den anderen Teilnehmer sicher ist.

...

## **E.3 Treatment: Force-Tradeoff**

...

### **Ihre Information**

Bevor Sie sich entscheiden, können Sie zusätzliche Information darüber erhalten, welche der Optionen der Computer für Sie als sicher ausgewählt hat.

...

## **E.4 Treatment: Force-Control**

...

## Ihre Entscheidungen

...

Ihre Entscheidung besteht daraus, zwischen zwei Optionen, X und Y, zu wählen.

...

4. Ihre eigene Auszahlungshöhe ist unabhängig davon, ob Sie sich für X oder Y entscheiden. Ihre Auszahlungshöhe ist auch unabhängig davon, ob Option X oder Option Y für den anderen Teilnehmer sicher ist.

...

## Ihre Information

Bevor Sie sich entscheiden, können Sie zusätzliche Information darüber erhalten, welche der Optionen der Computer für Sie als sicher ausgewählt hat.

...

## F Instructions English translation

### F.1 Treatment: NoForce-Tradeoff

### General Explanations

Welcome to the study! In this study, you can earn a good amount of money. Please carefully read the following explanations! Shall you had questions, please stick your hand out of the cubicle—we will come to your seat.

**During the study, it is not allowed to talk with other participant, to use mobile phones, nor to start other programs on the computer.**

The violation of these rules will lead to an exclusion form the study and any

payment. You will receive your payment of the study at the end of the study in cash. During the study, we do not talk about Euro. Instead we will talk about points. Your total payment will be calculated in points and translated into Euro at the following rate:

**1 Point = 5 Cent.**

**Participant Pairing:** The computer has paired you with another participant who is randomly selected from all the participants in the room. In the following, we refer to the participant whom you are paired with as ‘the other participant’.

Please feel free to refer back to this Instruction when you are working on the computer. Before the study starts, we will ask you to answer a couple of control questions.

## Probabilities

In this section, we want to familiarize you with mathematical probabilities.

Imagine the following situation: in a room there are 10 participants. Every participant has a box. The participants cannot see the boxes but know the following:

- In each box there are 10 balls.
- 5 of the 10 participants have a box with 10 white balls and 4 black balls (Situation A).
- 5 of the 10 participants have a box with 4 white balls and 6 black balls (Situation B).

Now imagine the following: you are one of the participants. So the probability that you are in Situation A is 50%. A computer draws a ball out of your box and places it back into the box after the draw. When you find out the color of the drawn ball, it helps to better assess what situation you are in.

Question: What is the probability that you are in Situation A, given that the drawn ball is white?

Probabilities, which take additional information into account—such as the color of the drawn ball—are calculated in a statistically correct manner according to a mathematical law, Bayes' theorem. The calculation is complicated and requires several calculation steps. **We therefore show you the statistically correct probabilities whenever they are relevant for your decision.**

Example (Answer to the Question):

$$\begin{aligned} & \text{Probability of Situation A, given that the drawn ball is white} \\ = & \text{Probability of Situation A} \cdot \left[ \frac{\text{Probability that the drawn ball in situation A is white}}{\text{Probability that the drawn ball is white}} \right] \\ = & 50\% \cdot \frac{60\%}{50\%} \\ = & 60\%. \end{aligned}$$

The probability of Situation A, given that the drawn ball is black, is calculated in a similar way. **Now, please process to the first control question on the computer.**

## Your Decisions

You obtain now 100 points to your points account.

Next you make a decision that affects your own payout amount and the payout amount of the other participant, who also received 100 points on the points account. However, the decisions of the other participant have no impact on your payout amount.

Your decision is to choose between two options, X and Y.

1. One of these options is 'harmful' to the other study participant and leads to 80 points less being paid out to the participant.
2. The other of the options has no effect on the other study participant,

this option is 'safe'.

3. For every 7 out of 20 decision-makers, option X is safe and option Y is harmful (35 % probability), and for every 13 out of 20 decision-makers, option Y is safe and option X is harmful (65 % probability ). The computer has already chosen at random which option is safe and which is harmful in your case.
4. Regardless of whether option X or option Y is safe for the other participant, you will receive 25 points more yourself if you choose option X.

## Your Information

Before you make a decision, you can get additional information about which of the options the computer has selected to be safe for you. If you want to make the decision without additional information, please click directly on 'Decision without additional information'. If you want to make the decision with additional information, please click on 'Decision with additional information'.

There are 100 white or black balls in a box. If option X is safe, there are 60 white balls and 40 black balls in the box. If option Y is safe, there are 40 white balls and 60 black balls in the box.

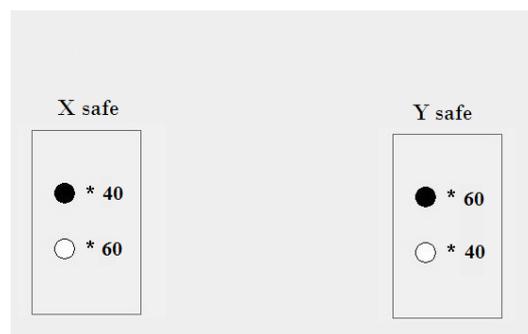


Figure 11: Box with 40 black and 60 white balls (X safe); Box with 60 black and 40 white balls (Y safe)

You won't be shown the box, but you can have the computer draw a ball

out of the box at random. To do this, click the button 'Another Ball'. After you click, the ball that the computer drew out of the box will appear on the screen. Then the computer puts the ball back in the box so that there are 100 balls in the box again.

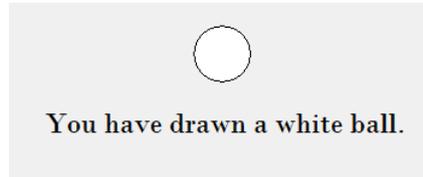


Figure 12: Example: The ball drawn by the computer is white.

You can then let another ball be drawn. To do this, click the 'Another Ball' button again. You can have an unlimited number of balls drawn. If you do not want to have any more balls drawn, click the button 'To the decision' .

Do you remember that there are more white balls in the box when option X is safe (60 white balls) than when option Y is safe (40 white balls)? So if you see that a white ball has been drawn, it is an indication that option X is safe. After each ball, the statistically correct probability that option X is safe is displayed. All balls that you have already drawn are taken into account.<sup>47</sup>

The person affected by your choice does not know whether or how much you have obtained information. You won't find out at the end of the study whether X or Y is the harmful action. It is free of charge for you to obtain additional information before making your decision between X and Y.

<sup>47</sup>A mathematical law, Bayes' theorem, is used to accurately calculate the probability that option X is safe. The exact calculation formula, if e.g. a single white ball has been drawn, is the following:

$$\begin{aligned}
 & \text{Prob. that option X is safe conditional on the ball drawn being white} \\
 = & \text{ Prob. that option X is safe} \\
 & \frac{\text{Prob. that a white ball is drawn conditional on option X being safe}}{\text{Prob. that a white ball is drawn}}.
 \end{aligned}$$

This is the only objectively correct calculation method.

## **F.2 Treatment: NoForce-Control**

### **Your Decisions**

...

Your decision is to choose between two options, X and Y.

...

4. Your own payout amount is independent of whether you choose X or Y. Your payout amount is also independent of whether option X or option Y is safe for the other participant.

## **F.3 Treatment: Force-Tradeoff**

...

### **Your Information**

Before you make a decision, you can get additional information about which of the options the computer has selected to be safe for you.

...

## **F.4 Treatment: Force-Control**

### **Your Decisions**

...

Your decision is to choose between two options, X and Y.

...

4. Your own payout amount is independent of whether you choose X or Y. Your payout amount is also independent of whether option X or option Y is safe for the other participant.

...

## **Your Information**

Before you make a decision, you can get additional information about which of the options the computer has selected to be safe for you.

...