

# DISCUSSION PAPER SERIES

DP17665

## **ANTICIPATORY ANXIETY AND WISHFUL THINKING**

Jan Engelmann, Maël LeBreton, Nahuel Salem-  
Garcia, Peter Schwardmann and Joël van der Weele

**ORGANIZATIONAL ECONOMICS**

**CEPR**

# ANTICIPATORY ANXIETY AND WISHFUL THINKING

*Jan Engelmann, Maël LeBreton, Nahuel Salem-Garcia, Peter Schwardmann and Joël van der Weele*

Discussion Paper DP17665  
Published 13 November 2022  
Submitted 31 October 2022

Centre for Economic Policy Research  
33 Great Sutton Street, London EC1V 0DX, UK  
Tel: +44 (0)20 7183 8801  
[www.cepr.org](http://www.cepr.org)

This Discussion Paper is issued under the auspices of the Centre's research programmes:

- Organizational Economics

Any opinions expressed here are those of the author(s) and not those of the Centre for Economic Policy Research. Research disseminated by CEPR may include views on policy, but the Centre itself takes no institutional policy positions.

The Centre for Economic Policy Research was established in 1983 as an educational charity, to promote independent analysis and public discussion of open economies and the relations among them. It is pluralist and non-partisan, bringing economic research to bear on the analysis of medium- and long-run policy questions.

These Discussion Papers often represent preliminary or incomplete work, circulated to encourage discussion and comment. Citation and use of such a paper should take account of its provisional character.

Copyright: Jan Engelmann, Maël LeBreton, Nahuel Salem-Garcia, Peter Schwardmann and Joël van der Weele

# ANTICIPATORY ANXIETY AND WISHFUL THINKING

## Abstract

We test the hypothesis that anxiety about adverse future outcomes leads to wishful thinking. Across four experiments (N=1,116), participants perform pattern recognition tasks in which some patterns may result in an electric shock or a monetary loss. Participants engage in significant wishful thinking, as they are less likely to correctly identify patterns that may lead to a shock or loss. Wishful thinking increases with greater ambiguity of the visual evidence and is only disciplined by higher accuracy incentives when accuracy depends on participants' cognitive effort. Wishful thinking is heterogeneous across and stable within individuals.

JEL Classification: N/A

Keywords: N/A

Jan Engelmann - [jbengelmann@gmail.com](mailto:jbengelmann@gmail.com)

*University of Amsterdam, University of Amsterdam & Tinbergen Institute*

Maël LeBreton - [mael.lebreton@googlemail.com](mailto:mael.lebreton@googlemail.com)

*Paris School Of Economics*

Nahuel Salem-Garcia - [salemnahuel@gmail.com](mailto:salemnahuel@gmail.com)

*University of Geneva*

Peter Schwardmann - [pschwardmann@gmail.com](mailto:pschwardmann@gmail.com)

*Carnegie Mellon University and CEPR*

Joël van der Weele - [j.j.vanderweele@uva.nl](mailto:j.j.vanderweele@uva.nl)

*University of Amsterdam, University of Amsterdam & Tinbergen Institute*

## Acknowledgements

The authors thank the editor, four anonymous referees, Douglas Bernheim, Andrew Caplin, Mark Dean, Yves Le Yaouanq, George Loewenstein, Nathaniel Neligh, Matthew Rabin, Simeon Schudy, Claudia Senik, Severine Toussaert, and seminar participants at the 2021 AEA Meeting, Stanford University, Amsterdam Brain and Cognition Center, ECBE, University of Chicago, New York University, MBEES workshop, MiddExLab Seminar, Belief Based Utility Conferences at BRIQ and UvA, Paris School of Economics, LMU Munich, and UC Santa Barbara for useful comments. Li-Ang Chang, Sara Khayouti and Nik Mautner Markhof provided excellent research assistance. For financial support, Peter Schwardmann thanks the Bavarian Academy of Sciences and Humanities, Joel van der Weele thanks the Dutch Science Association (NWO) via a personal VIDI grant (452-17-004), and Mael Lebreton thanks the Swiss National Science Foundation (SNSF) via a personal Ambizione grant (PZ00P3\_174127).

# Anticipatory Anxiety and Wishful Thinking\*

Jan B. Engelmann<sup>1</sup>, Maël Lebreton<sup>2</sup>, Nahuel A. Salem-Garcia<sup>3</sup>, Peter Schwardmann<sup>4</sup>, and  
Joël J. van der Weele<sup>1</sup>

<sup>1</sup>University of Amsterdam, Tinbergen Institute

<sup>2</sup>Paris School of Economics

<sup>3</sup>University of Geneva

<sup>4</sup>Carnegie Mellon University

October 31, 2022

## Abstract

We test the hypothesis that anxiety about adverse future outcomes leads to wishful thinking. Across four experiments (N=1,116), participants perform pattern recognition tasks in which some patterns may result in an electric shock or a monetary loss. Participants engage in significant wishful thinking, as they are less likely to correctly identify patterns that may lead to a shock or loss. Wishful thinking increases with greater ambiguity of the visual evidence and is only disciplined by higher accuracy incentives when accuracy depends on participants' cognitive effort. Wishful thinking is heterogeneous across and stable within individuals.

**JEL classification:** C91, D83

**Keywords:** confidence, beliefs, anticipatory utility, anxiety, motivated cognition

---

\*The authors thank the editor, four anonymous referees, Douglas Bernheim, Andrew Caplin, Mark Dean, Yves Le Yaouanq, George Loewenstein, Nathaniel Neligh, Matthew Rabin, Simeon Schudy, Claudia Senik, Severine Toussaert, and seminar participants at the 2021 AEA Meeting, Stanford University, Amsterdam Brain and Cognition Center, ECBE, University of Chicago, New York University, MBEES workshop, MiddExLab Seminar, Belief Based Utility Conferences at BRIQ and UvA, Paris School of Economics, LMU Munich, and UC Santa Barbara for useful comments. Li-Ang Chang, Sára Khayouti and Nik Mautner Markhof provided excellent research assistance. For financial support, Peter Schwardmann thanks the Bavarian Academy of Sciences and Humanities, Joël van der Weele thanks the Dutch Science Association (NWO) via a personal VIDI grant (452-17-004), and Maël Lebreton thanks the Swiss National Science Foundation (SNSF) via a personal Ambizione grant (PZ00P3<sub>1</sub>74127).

# 1 Introduction

Many commonly held beliefs seem to be inspired by their comforting properties rather than their realism. Billions of adherents of the major religions believe in an afterlife, despite a lack of evidence for its existence. Moreover, religiosity is higher in populations that face unpredictable shocks like earthquakes (Sinding Bentzen, 2019), during pandemics (Sinding Bentzen, 2021), and in the absence of alternative forms of insurance (Auriol et al., 2017). People at risk of serious diseases avoid medical testing and remain optimistic about their health status (Lerman et al., 1998; Oster et al., 2013; Ganguly and Tasoff, 2016), while greater exposure to Covid-19 leads people to become more sanguine about the probability of infection (Orhun et al., 2021; Islam, 2021). Populist politicians that promise easy fixes find more support in areas with weak economic prospects and declining growth rates (Mughan et al., 2003; Obschonka et al., 2018).

These findings suggest that the adoption of comforting beliefs affects important decisions and originates in anxiety about adverse future outcomes. At the same time, it is hard to establish with field data that beliefs are the product of self-deception or wishful thinking, or to determine the precise motives behind self-deception.<sup>1</sup> Laboratory experiments on wishful thinking or optimism bias have also not established a clear causal link with anxiety, due to a number of factors that we discuss in the next section. This leaves an important gap in the literature: pinning down the motives and processes behind wishful thinking allows predicting its occurrence as well as the design of strategies that may abate it, such as tackling the underlying anxiety, providing precise information, or raising the material costs of false beliefs.

To make progress on these matters we study the effect of experimentally induced anxiety on belief formation in a tightly controlled setting. Our four preregistered experiments (combined  $N = 1,116$ ) incentivize participants to correctly identify which of two types of patterns they see on their screen. We induce anxiety by associating one type of pattern with an adverse outcome that may occur after a short waiting period. In our first experiment, the adverse outcome is a mild electric shock, a proven method of inducing anxiety. In our second, third, and fourth experiment, the adverse outcome is a monetary loss. Since participants have no control over the occurrence of these adverse outcomes, the payoff-maximizing strategy is to identify the patterns as accurately as possible. By contrast, anticipatory anxiety about the shock or loss may cause wishful thinking, a

---

<sup>1</sup>For instance, consistent with wishful thinking, Oster et al. (2013) find that people at risk of Huntington disease are optimistic before they get tested for the disease, but are reluctant to test, especially when they have low objective risk. However, without exogenous variation in the motives to hold optimistic beliefs, it is neither clear whether initial optimism is the result of wishful thinking nor whether it is driven by a desire to avoid feeling anxious. Furthermore, Islam (2021) finds that individuals self-deceive about the risk of a Covid infection in deciding whether to go to a coffee shop during the pandemic. At the same time, they distort beliefs about the risk for others rather than for themselves, suggesting that self-deception is driven by social motives rather than anxiety about one's own health.

belief that the anxiety-inducing state of the world is less likely than it really is. In our experiment, wishful thinkers will then be less accurate when the pattern that is flashed on the screen is associated with a shock or monetary loss, and more accurate when the pattern that is not flashed is associated with a shock or loss.

We propose a simple model to clarify the properties of wishful thinking in our experimental setting. Following Bénabou and Tirole (2002) and Brunnermeier and Parker (2005), we suppose that an agent distorts the visual signal to optimally trade off the anticipatory utility benefits from alleviated anxiety and the material costs stemming from inferior decision making. Apart from wishful thinking or overoptimism, this generates a number of additional predictions, most notably about the role of ambiguity and accuracy incentives, which are shared by most models motivated cognition. We experimentally test these predictions across our experiments, and develop an individual measure of wishful thinking of our participants.

The results show clear evidence for wishful thinking. In all experiments, participants are significantly less accurate in identifying patterns that may lead to an adverse outcome. This result obtains for different sources of anxiety (shock versus monetary loss), different pattern identification tasks, and in different settings (the experiment taking place online versus in the laboratory). Our dataset can rule out competing explanations for the observed effect, like illusions of control, whereby participants believe that the pattern they report determines the adverse outcome, or the idea that adverse outcomes scare participants into providing noisy responses. We also find that, on balance, wishful thinking does not depreciate and remains high in later trials of the experiments. These results are strikingly robust in comparison to the literature in economics and psychology that we review below, which has yielded mixed findings, mostly in the domain of positive outcomes.

In all experiments, we find that wishful thinking is more pronounced when the evidence is more ambiguous - i.e. when the different patterns are difficult to distinguish. This is true across three different visual inference tasks in which we vary ambiguity in distinct ways. Therefore, in line with previous results (Haisley and Weber, 2010; Sloman et al., 2010) and recent theoretical work (Caplin and Leahy, 2019), we find that the precision of the signal constrains motivated belief formation.

We also test whether wishful thinking decreases with increasing material costs of false beliefs. We vary the accuracy bonus that participants can earn from a correct answer by factors up to 200. In our first three experiments, higher accuracy incentives do not lead to a decrease in wishful thinking. They also do not lead to an increase in accuracy. This is not for lack of trying on behalf of the participants, as response times and self-reported concentration increase significantly under higher incentives. In Experiment 2, we also vary the magnitude of monetary losses, i.e. the adverse outcome that participants may feel anxious about or the "anticipatory payoffs". Similar to our

results on material incentives, we find that increasing losses increases self-reported anxiety but has no significant effect on wishful thinking.

We interpret these null results in the context of our model, where the agent first observes a signal and then distorts her mental representation of it to become more optimistic. The results of our first three experiments imply that participants engage in this ex-post signal distortion, but that the extent of the distortion does not respond to either material or anticipatory payoffs at the margin. However, our results suggest another way in which wishful thinking may respond to accuracy incentives. Incentives increase participants' ex-ante effort to form an accurate representation. If these efforts are successful, they may make wishful thinking harder in the same way that being shown a less ambiguous pattern increases accuracy and reduces the scope for wishful thinking in our experiments. In our first three experiments, this mechanism is ruled out by the fact that accuracy is insensitive to cognitive effort.

To test this alternative mechanism, we designed Experiment 4 to make accuracy maximally elastic in effort: the pattern recognition task is self-timed and participants can arrive at a correct answer through a laborious counting exercise. Here, we find evidence that incentives reduce wishful thinking, precisely when incentives lead to increased accuracy. In particular, when we focus on about 40 percent of participants who count more often in response to higher accuracy incentives, we find that these participants, but not others, actually increase their average accuracy and reduce their wishful thinking. These results suggest that material payoffs affect wishful thinking only when they lead to investment and subsequent improvement in signal precision, but not by affecting the inner calculus of ex-post signal distortion. This insight adds nuance to a small but growing literature on the role of accuracy incentives in disciplining motivated beliefs, which, as we discuss in the next section, has yielded several null results.

A final set of results concerns wishful thinking as a personal characteristic or trait. Because participants go through many trials, we can compute individual-level measures of wishful thinking to study heterogeneity in people's proclivity to engage in motivated cognition, a novelty in the experimental literature on this topic. We find that wishful thinking is stable within individuals, but variable across them, with some individuals displaying the opposite tendency. Moreover, an individual's wishful thinking correlates negatively with their self-reported concentration on the pattern recognition task.<sup>2</sup> This is consistent with the idea that ex-ante investments in signal precision decrease wishful thinking. Another intuitive finding is that wishful thinking correlates negatively with a measure of defensive pessimism, a tendency to expect the worst in order to

---

<sup>2</sup>Therefore, the "natural" differences in concentration between individuals appear to be more consequential than the within-subject shifts in concentration that a higher accuracy bonus managed to induce in the first three experiments.

avoid disappointment. Instead, the correlation between experienced anxiety and wishful thinking is positive. We also find suggestive evidence for a positive correlation between wishful thinking and risk seeking as well as between wishful thinking and the belief in an afterlife, but no correlation with worries about climate change.

In the next section, we review the literature on optimism and wishful thinking. We then describe our experimental design. Section 4 introduces a simple theoretical model that helps us derive our hypotheses. Section 5 contains the results of our experiment before section 6 investigates heterogeneity in wishful thinking. We conclude in section 7.

## 2 Literature

The literature on wishful thinking, optimism bias, and desirability bias has yielded mixed results and there is still an active debate about the phenomenon’s existence, scope and its underlying mechanisms. One strand of this literature uncovers optimism in probability judgements about real life events (Weinstein, 1982; Lench and Ditto, 2008). Other studies claim behavioral and neural evidence for asymmetric updating about future life events, whereby bad news is downweighted (Sharot et al., 2011, 2012). However, these results have been called into question, with critics suggesting that they can be explained by standard Bayesian updating (Shah et al., 2016; Burton et al., 2022). Economic experiments have used more stylized information structures that can rule out Bayesian updating patterns Möbius et al. (2014). However, experiments that study updating from ego-relevant information, such as scores on an IQ test, have yielded mixed results (see Drobner (2022) for a review).

Psychologists have studied optimism using the “marked-card” paradigm, wherein participants rate the probability of drawing a particular card that is associated with monetary outcomes. In a meta-analysis, (Krizan and Windschitl, 2007) find evidence for optimism or desirability bias in this paradigm, but not in other, related settings (see also Windschitl et al., 2010). Some papers on “motivated perception” induce biased perceptions of ambiguous visual evidence (e.g., an image that could be interpreted as a B or a 13) by telling participants that one interpretation of the evidence results in the consumption of a preferred drink or food (Balcetis and Dunning, 2006). These studies cannot incentivize beliefs because there is no true state of the world and instead rely on implicit questionnaire items, eye tracking and reaction times (Dunning and Balcetis, 2013) to get at deeply held beliefs. These studies also struggle to rule out that participants believe that their answers can affect outcomes. Leong et al. (2019) shows that monetary prizes affect visual perceptions and provides neurological evidence about the location of the perceptual distortions in the brain.



In economics there is a small number of studies on wishful thinking in the domain of monetary gains, again yielding mixed results. Mijović-Prelec and Prelec (2010), Coutts (2019), and Mayraz (2011) document overoptimism in estimation tasks in which participants have an exogenous monetary stake in some of the outcomes. However, Coutts (2019) finds it only for one out of three tasks. Barron (2021) finds no evidence for asymmetries in updating of beliefs about the probability of winning monetary prizes.

Our paper differs from previous laboratory studies in several ways. Most importantly, we focus on anxiety as a driver of wishful thinking and on negative rather than positive outcomes. The anticipation of negative events may have a more powerful influence on belief formation both because it may activate different cognitive processes than the anticipation of gains and because people tend to care more about losses than about equivalent gains.<sup>3</sup> This may explain why wishful thinking emerges as a robust phenomenon across our experiments, tasks and contexts, which contrasts sharply with the mixed results in the literature.

Our design also introduces a number of other innovations. First, electric shocks are a proven way of inducing anxiety and allow precise control over its timing. Second, we vary the ambiguity of evidence in a subtle and inconspicuous way. Third, we administer within-subject treatments with many observations per person, which allows us to look at wishful thinking as an individual characteristic that may be correlated with other traits. Earlier work in related settings is scarce and has not yielded consistent results. Buser et al. (2018) do not find significant correlations between asymmetric updating of ego-relevant news across three tasks. However, there are few observations per participant and the repeated updating task is noisy and subject to other biases like conservatism and base rate neglect. Sharot et al. (2007) and (Sharot et al., 2012) show neural and hormonal substrates of optimism bias, suggesting a hardwired component that may differ between individuals.

Like us, previous experiments have investigated whether accuracy incentives reduce motivated beliefs, in line with a trade-off between the psychological motive for and the (material) costs of adopting wrong beliefs that is central to models of motivated beliefs (e.g. Bénabou and Tirole 2002; Brunnermeier and Parker 2005; Bénabou and Tirole 2011). In the affirmative, Armor and Sackett (2006) find more optimism for hypothetical than real events and Zimmermann (2020) shows that incentives can reduce motivated biases in recall. However, much evidence goes in the other direction. Simmons and Massey (2012) show that accuracy incentives of up to \$50 do not correct football fans' overoptimistic expectations about their home team. Lench and Ditto (2008) find

---

<sup>3</sup>Falk and Zimmermann (2016) study the role of anxiety in information preferences and investigate whether participants have a preference over early or late resolution of uncertainty about the occurrence of an electric shock.

no effect of incentives on optimistic beliefs about adverse life events. Mayraz (2011) and Coutts (2019) find that higher rewards for accuracy do not reduce wishful thinking, and Schwardmann et al. (2022) finds no evidence for an effect on self-persuasion and polarization in a debating context.<sup>4</sup>

Our results provide a new and more nuanced view of the role of incentives: We find that accuracy incentives only reduce motivated beliefs in tasks where participants can improve the precision of signals through effort and thereby reduce the scope for wishful thinking. This suggests that the impact of economic incentives on motivated beliefs is likely to be highly sensitive to the nature of the inference task and the extent to which accuracy is elastic in effort.

### 3 Design

Our study comprises four computerized experiments, which we number in the order in which they were run. We preregistered hypotheses for each experiment on [Aspredicted.org](https://aspredicted.org). The preregistrations and IRB approvals can be found in Appendix C.4.2. All experimental instructions are available in Appendix E.

#### 3.1 Design features common to all experiments

All experiments share the same basic structure. In each experiment, participants engaged in a number of trials of a pattern recognition task. In each trial, they had to identify which of two possible types of pattern was shown on the screen. One of the two patterns was associated with an undesirable outcome: an electric shock or a monetary loss, depending on the experiment. We refer to trials in which the pattern associated with a shock or loss and the pattern that was flashed on the screen were aligned as “shock/loss patterns” and trials in which they were not aligned as “no-shock/no-loss patterns”.

If the no-shock/no-loss pattern was shown, then no shock or loss would occur in the trial. If a shock/loss pattern was shown on the screen, then the shock or loss occurred with a probability of one third at any point within an eight second period following the participants’ response to the trial. This procedure injects objective uncertainty into the occurrence of the shock or loss. The probabilistic implementation also assures that shocks occur sparingly, which avoids rapid desensitization (or sensitization) of participants. Because participants will generally not be completely certain which pattern they saw, there is additional subjective uncertainty. In keeping with the pre-

---

<sup>4</sup>More generally, (Enke et al., 2021) investigate a number of well-known cognitive biases and show that paying up to monthly salary for accuracy does not improve performance, although it does induce more cognitive effort. This confirms an earlier review that concluded that “no replicated study has made rationality violations disappear purely by raising incentives” (Camerer and Hogarth, 1999).

vious literature, we will refer to the emotions the uncertain shock or loss induces as “anticipatory anxiety”.<sup>5</sup>

Our first and main treatment varies the associations between patterns and shocks or losses. Between trials and within participants, we varied not just the actual pattern but also which type of pattern was associated with a shock or loss. This assured that any differential response to the two types of patterns could not affect our results. The occurrence of the shock depended only on the pre-determined shock pattern and the actual pattern on the screen and not on the participant’s response. However, believing that one saw a no-shock pattern could reduce anxiety about the imminent shock or loss. This would bias participants answers towards no-shock patterns, and hence make them less accurate when a shock pattern is shown and more accurate when a no-shock pattern is shown. We measure wishful thinking as the difference between average accuracy for “no-shock” and “shock” patterns.

Each experiment featured at least two further within-subject treatment variations. One of these varied the ambiguity of the pattern, in order to test whether wishful thinking is stronger for more difficult/ambiguous patterns. Another treatment varied the bonus that participants could win for a correct response, resulting in a *Low Accuracy Bonus* and a *High Accuracy Bonus* condition. This experimentally manipulated the trade-off between psychological payoffs from having more optimistic beliefs and the material payoffs from having more accurate beliefs. The order of these treatments was fully counterbalanced in each experiment. Participants received no explicit feedback about their performance.

Each experiment also implemented a series of variations on this basic structure in order to answer specific research questions. We summarize these variations in Table 1 and discuss them in turn.

### 3.2 Experiment 1: Electric Shocks

The experiment took place in the CREED experimental laboratory at the University of Amsterdam. Sixty subjects participated in individual sessions. Upon coming to the lab, subjects read the instructions, signed a consent form and answered several control questions to determine their

---

<sup>5</sup>The American Psychological Association defines anxiety as “worry or apprehension about an upcoming event or situation because of the possibility of a negative outcome, such as danger, misfortune, or adverse judgment by others.” The clinical psychology literature sometimes makes a distinction between fear and anxiety. Fear is defined as a behavioral response that serves to mobilize an organism in life-threatening situations that present immediate and identifiable danger. Anxiety, on the other hand, produces a more sustained response to aversive events that are unpredictable in terms of their timing and frequency, resulting in prolonged worry, tension and a feeling of insecurity (Grillon, 2008; Schmitz and Grillon, 2012). However, the fine points of the distinction differ between authors, and threats may induce a mixture of these emotions. Indeed, our design implements some elements of fear-induction (the threat is a clearly identifiable shock or loss) and anxiety-induction (the shock or loss is uncertain).

Table 1: Overview of experimental designs

	Experiment 1	Experiment 2	Experiment 3	Experiment 4
<b>Participants</b>	60	221	426	409
<b>Number of trials</b>	216	up to 96	up to 64	up to 96
<b>Visual task</b>	Single Gabor flash	Single Gabor flash	8 Gabor flashes	Colored dots
<b>Anxiety source</b>	Electric shock	Monetary loss	Monetary loss	Monetary loss
<b>Loss/shock size</b>	Self-calibrated	0, 0.1 or 5 pounds	1 pound	0 or 1 pound
<b>Task difficulty levels</b>	Tilt size (3 levels)	Tilt size (2 levels)	Likelihood ratio (continuous)	Dot ratio (4 levels)
<b>Accuracy bonus levels</b>	1 euro 20 euro	10 cents 8 pounds	5 cents 10 pounds	5 cents 10 pounds
<b>Other design elements</b>	Confidence measure Replication exp.		Treatment reminders	Treatment reminders Self-timed task
<b>Start / end date</b>	November 12, 2018 December 5, 2018	Feb 23, 2021 Feb 29, 2021	Jan 3, 2022 Jan 4, 2022	March 8, 2022 March 8, 2022
<b>Location</b>	Laboratory (Amsterdam)	Online (Prolific.co)	Online (Prolific.co)	Online (Prolific.co)

understanding of the task and the belief elicitation mechanism. The experimenter pointed out any wrong answers and discussed the correct answer until the participant indicated they understood them.

The source of anxiety in this experiment was a mild electric shock. Electric shocks are a proven method of inducing anticipatory anxiety (Grillon, 2008; Schmitz and Grillon, 2012; Engelmann et al., 2015, 2019).<sup>6</sup> Moreover, they are salient consumption events that afford a lot of control over the precise timing of the emotions. Since people differ in their pain thresholds, the strength of the electric shock was calibrated individually.<sup>7</sup>

The visual task was to determine whether a grating (Gabor patch), was tilted towards the left or right (see example in Panel (a) of Figure 1). Before each trial, subjects were reminded of the treatment conditions. After briefly seeing a fixation cross (750ms), the grating was flashed on the

<sup>6</sup>In particular, people pay to shorten the time they have to wait for electric shocks (Loewenstein, 1987; Berns et al., 2006) and they display physiological arousal while waiting for them, as reflected in a heightened skin conductance response (Engelmann et al., 2015, 2019).

<sup>7</sup>The wrist of the participant’s non-dominant hand was connected to a Digitimer DS5 isolated bipolar current stimulator, which itself was connected to MATLAB through National Instruments USB x-series. The participant induced herself with a series of shocks, which she rated on a pain scale of 0 (not painful at all) to 10 (extremely painful). The calibration was complete when the subject rated the pain as 7-9 on the scale three consecutive times. A rating of 10 would lead to a decrease in the threshold. The maximum possible shock strength was set to 5V 25mA and the duration of the shock was set to 50ms (Engelmann et al., 2015, 2019).

screen (150ms). Participants were then asked to indicate the direction of the tilt by pressing the left or right arrow on the keyboard (self-paced) as well as the confidence in their choice on a scale from 50% (completely uncertain) to 100% (certainty). We incentivized confidence ratings with a Becker-deGroot-Marschak (BDM) or “matching probabilities” mechanism. This mechanism makes it incentive compatible to state true beliefs, regardless of a participant’s risk preferences.<sup>8</sup>

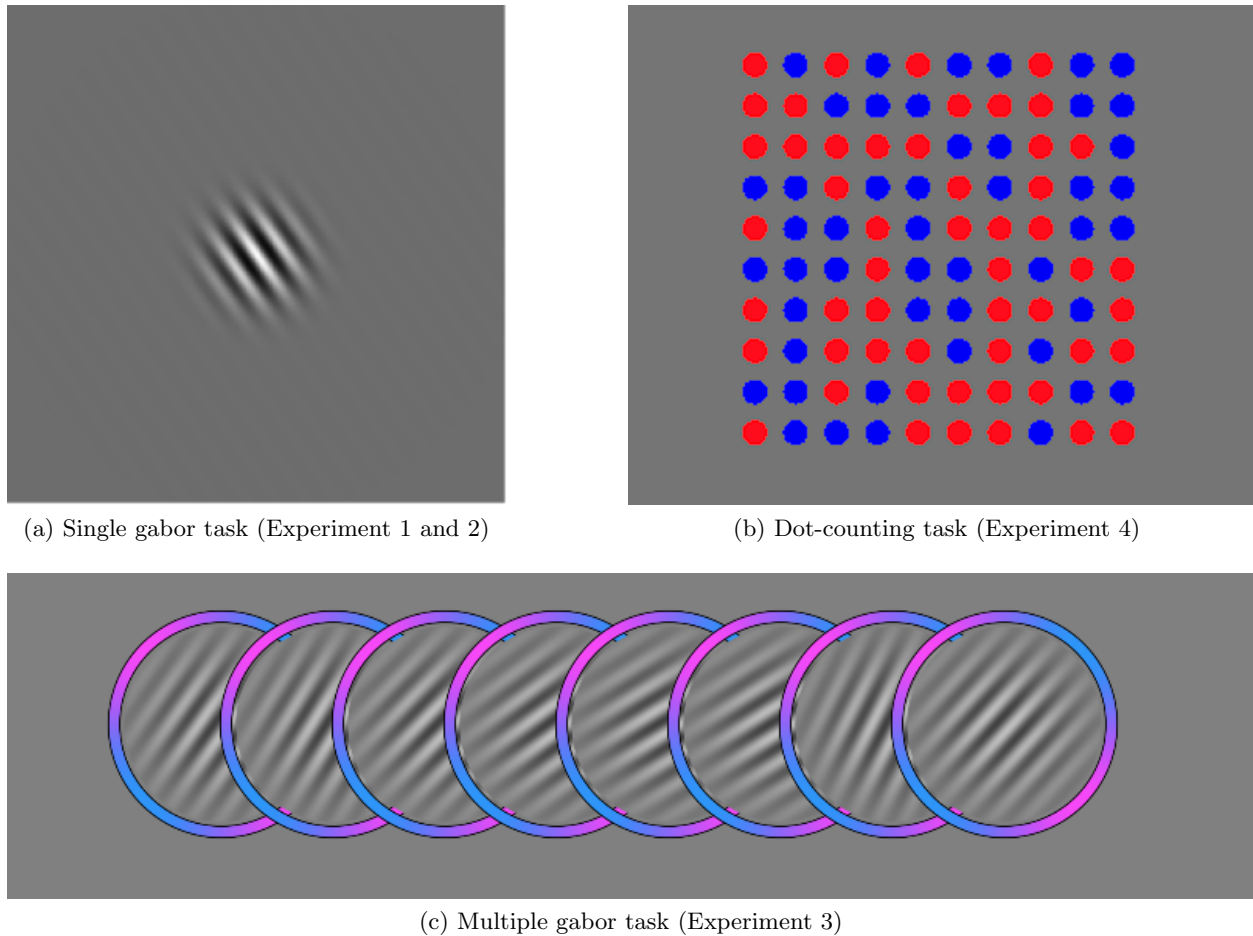


Figure 1: Examples of the visual tasks in the various experiments.

Next, participants faced an anticipation screen (2000-8000ms), asking them to wait for the shock resolution. Finally, the electric shock was administered or not (1000ms). No trial-by-trial feedback was given about the correctness of the guess, but the average performance was communicated at

---

<sup>8</sup>Subjects indicate their subjective probability  $x \in \{50, 55, \dots, 95, 100\}$  that their answer was correct. The computer then randomly draws a number  $z \in [50, 100]$ . If  $x \geq z$ , then subjects win prize  $M$ , if their answer truly is correct. If  $x < z$ , then subjects win prize  $M$  with probability  $z$ .  $M$  varies between experimental conditions. Schlag et al. (2015) provide details about the origins and incentive compatibility of this mechanism. See Trautmann and van de Kuilen (2014) and Hollard et al. (2016) for evidence. After the instructions but before the experiment started, participants had the opportunity to gain experience with the BDM mechanism.

the end of each block of 18 trials. Each participant completed three sessions, each divided in four blocks of 18 trials. The four blocks correspond to four conditions of a 2x2 factorial design (Shock x Incentive). As described above, the Shock treatment varied whether the possibility of a shock was associated with a right-tilted or left-tilted grating pattern. The Incentive treatment varied whether the potential prize in the belief elicitation was 1 or 20 euros. We also varied the difficulty of the pattern recognition task within each block, as measured by the degree of the tilt from the vertical line.<sup>9</sup>

Participants' earnings consisted of a 10 euro show-up fee, plus the earnings from the accuracy payments of one randomly drawn trial from both the low and high incentive condition. Thus, payments varied between 10 and 31 euros for a session that lasted on average slightly over an hour.

### 3.3 Experiment 2: Monetary losses as a source of anxiety

While electric shocks are a proven way to induce anxiety, they are not a common occurrence in everyday life. It is therefore important to understand whether the phenomenon carries over to other sources of anxiety, for instance the prospect of monetary losses.<sup>10</sup> Experiment 2 investigates wishful thinking in the presence of monetary losses. The experiment took place online, with 221 participants recruited from the online platform Prolific, which assures the highest quality of online data provision (Eyal et al., 2021). Participants had to answer a number of attention checks to advance to instructions, and a number of quiz questions about the instructions to advance to the experiment (see Appendix E). All monetary amounts were communicated in pounds.

To implement losses, participants were endowed with an amount of money, and on each trial they could lose part of this endowment. Participants were confronted with the same Gabor visual task as in Experiment 1. If a "loss pattern" appeared on the screen, then the participant would lose 20% of the endowment with a probability of one third. As before, subjects had to wait up to 8 seconds to learn whether they lost the money. To make losses salient, they were accompanied by an animation of an exploding bag of money. The experiment was divided into three parts of up to 32 trials. If the participant ran out of endowment before the 32 trials, then the remaining trials were cancelled.

---

<sup>9</sup>The three difficulty levels were calibrated to result in accuracy levels of 60%, 70% and 80%. Initially, these levels were calibrated on the basis of a pilot, and were the same for all subjects. To reduce the effects of fatigue or learnings, difficulty levels were recalibrated for each subject after each part, using a logistical performance function. This happened without subjects' knowledge, so this aspect of the design could not be gamed. We dropped the (re)calibration in the other experiments.

<sup>10</sup>As we discuss in Section 2, the connection between monetary outcomes and optimism has previously been investigated by other papers in the positive domain, e.g. Mayraz (2011); Barron (2016); Coutts (2019), which has led to mixed findings.

Using monetary losses allows us to vary the size of the losses, and possibly the anxiety associated with these losses. To this end, each participant went through three parts of the experiment that varied in endowment and loss size: 25 pound endowment with 5 pound losses (*High Loss* condition), a 50 cents endowment with 10 cents losses (*Low Loss* condition), and no endowment with no threat of losses (*Neutral* condition). The part without losses served to address potential confounds that we discuss in Section 5.5.3.

To vary task difficulty, we used two different angles for the tilt of the pattern (with the steeper one being closer to the vertical line and hence harder to recognize). The accuracy incentives varied between trials to be either 8 pounds sterling or 10 cents. Unlike in the previous experiment, we did not elicit confidence measures. Instead, we randomly selected one 8 pound trial and one 10 cent trial and paid subjects if their answer was correct. We made this change to implement the most parsimonious design that still allows for our various treatment dimensions, while avoiding attrition, fatigue and confusion of online participants due to the time-consuming and involved instructions of the confidence elicitation.

All treatments, including the three parts with different endowment sizes, were administered within-subject in randomized order. In order to reduce cognitive load, the tilt of the loss pattern (left vs. right) and the incentive for accuracy were varied at the block level, where a block consisted of 8 trials. At the start of each block, subjects were informed of the loss tilt, accuracy incentives and loss size, and were shown a reminder before the start of each individual trial. At the end of each block we conducted an interblock survey in which we asked participants for their agreement with two statements, measured on a five point Likert scale. The first stated that subjects were anxious to lose money from their endowment, the second that they were concentrated on the task.

### **3.4 Experiment 3 and 4: Task characteristics and incentive effects.**

Next to the source of anxiety, a second dimension of robustness concerns the visual decision-making task. The nature of the task matters for two reasons. First, if we are to take wishful thinking seriously as a cognitive phenomenon, we should ascertain that it is robust across multiple tasks, in contrast to evidence in Coutts (2019). Second, the task may affect mental trade-offs and hence the effect of accuracy incentives. In particular, incentives may reduce bias by motivating people to work harder to obtain evidence and thereby increase their accuracy, which then reduces their capacity for wishful thinking. Our quickly flashed Gabor pattern may not provide enough levers for increasing performance and may therefore not provide a good test of this mechanism.

To investigate these issues in more detail, we selected two new tasks that draw on more effortful cognitive processes. To better test the effect of accuracy incentives, we reduced potential

distractions in treatment variation by keeping loss sizes fixed. We also highlighted the accuracy incentive variation by alerting subjects explicitly that performance on high bonus trials was more lucrative.<sup>11</sup>

**Experiment 3: Memory and inference task.** The task in Experiment 3 is based on Drugowitsch et al. (2016) – see also Salvador et al. (2022). Participants saw a consecutive sequence of 8 Gabor patches spaced over 4 seconds, as illustrated in Figure 1. The tilts were generated from one of two distributions of patterns that was biased towards either left or right-leaning patterns. We then asked participants to infer which distribution generated the patterns, and define a correct answer as the one that corresponds to the distribution with the highest posterior likelihood given the displayed patterns.<sup>12</sup>

This task requires memorizing and mentally combining the several cues, which has been identified as a bottleneck of decision accuracy beyond the visual processing and choice implementation steps that were the focus of our previous task (Drugowitsch et al., 2016; Findling and Wyart, 2021; Wyart and Koechlin, 2016). It therefore requires a new dimension of mental effort, through which incentives for accuracy may increase decision accuracy and/or reduce bias. This design builds on evidence that incentive effects are larger for more complex tasks (Garbers and Konradt, 2014).

The design of losses followed that of Experiment 2. Participants completed two parts. In each part they received an endowment of 5 pounds from which they would lose 1 pound with a probability of one third if a “loss pattern” appeared. The part finished when the endowment was exhausted (after 5 losses) or after 32 trials. Within each part of the experiment, there were up to four 8-trial blocks across which we varied the size of the accuracy bonus (0.05 cents vs. 10 pounds) and the orientation of the loss patterns (left vs. right). After each block, there was an interblock survey about concentration on the task. We recruited 426 subjects on Prolific, using the same procedures as in Experiment 2.

**Experiment 4: Dot task.** To further increase the link between mental effort and performance, we introduce a dot-counting task, displayed in Figure 1. Participants saw an array of 100 dots and were asked to identify whether the majority of dots were blue or red. The task is self-timed, with a time limit of 40 seconds. This allows participants to exercise a lot of control over their performance through the time they spend on verifying the correct answer, including by counting the dots on

---

<sup>11</sup>Instructions mentioned that “High Prize trials have a stronger impact on earnings than Low Prize trials. Participants who focus more on High Prize trials earn more on average than those who focus more on Low Prize trials.”

<sup>12</sup>Occasionally, this might differ from the actual distribution that generated the pattern, but in contrast to Drugowitsch et al. (2016), we focus on the correct answer from the perspective of the participant.



the screen. Perhaps for that reason, previous studies using these or very similar tasks have found strong effects of incentives for accuracy (Caplin and Dean, 2014; Dean and Neligh, 2019; Dewan and Neligh, 2020). In addition, Bosch-Rosa et al. (2021) find that self-image concerns lead to motivated belief formation in this task.

Other parts of the design followed that of Experiment 2 and 3. Participants completed two parts with 32 trials each. In each part, participants received an endowment of 5 euros from which they would lose 1 euro with a probability of one third if a “loss pattern” appeared. When a subject exhausted the endowment (i.e. after 5 losses), the part stopped. Within each part of the experiment, there were up to four 8-trial blocks across which we varied the size of the accuracy bonus (0.05 cents vs. 10 pounds) and the color of the loss pattern (blue vs. red). We varied the difficulty of the task, by varying whether the majority color has 51, 52, 53 or 54 dots. In addition, we included one “neutral” part of 32 trials without endowments or losses, the order of which was randomized to be either before or after the two parts with loss trials. Experiment 4 also featured the intertrial self-reports about anxiety and concentration that we used in Experiment 2. For experiment 4, we recruited 409 participants on Prolific.

## 4 Theoretical predictions

In this section, we present a stylized model of wishful thinking that captures our laboratory context and allows us to derive our main hypotheses. We will focus on the setting of experiment 1 and suppose that the threat of electric shocks is the source of anxiety. The model is in the spirit of Brunnermeier and Parker (2005). The agent chooses her beliefs, trading of the anticipatory utility benefits of optimism with the material costs stemming from wrong decisions. Moreover, belief distortions come at a cognitive cost as in Bénabou and Tirole (2002) and Bracha and Brown (2012).

The state of the world is given by  $r_\theta \in \{0, 1\}$ , where  $r_\theta = 1$  means that the true pattern is *right-tilted*. A participant observes a pattern and forms an initial probabilistic belief that  $r_\theta = 1$ , which we denote by  $p(r_\theta, s) \in [0, 1]$ . These undistorted initial beliefs  $p(r_\theta, s)$  depend on the true state  $r_\theta$ , with  $p(r_\theta = 1, s) \geq 0.5$  and  $p(r_\theta = 0, s) \leq 0.5$ . They also depend on the precision of the visual signal  $s$ , with  $\frac{dp(r_\theta=1,s)}{ds} > 0$  and  $\frac{dp(r_\theta=0,s)}{ds} < 0$ . In particular, they become more certain when the signal is more precise.

After perceiving the pattern and forming her initial beliefs, the agent self-deceives into a new belief  $\hat{p} \in [0, 1]$ . Assuming that the agent states her chosen belief  $\hat{p}$ , the Becker-DeGroot-Marshak (BDM) mechanism implies the following expected material payoffs from potentially winning a prize

$M$

$$\pi(p, \hat{p}) = \frac{1}{2} (1 + 2p\hat{p} - \hat{p}^2) M$$

The probability of winning the prize is maximized at  $\hat{p} = p$ . Therefore, if material payoffs were the only object in the agent's utility function, then she would not self-deceive.<sup>13</sup>

The agent's anxiety of the electric shock is based only on her chosen beliefs  $\hat{p}$  and is given by

$$\sigma_z(r_z\hat{p} + (1 - r_z)(1 - \hat{p}))qZ$$

The parameter  $\sigma_z \geq 0$  captures the importance of anticipatory utility concerns, or a participant's innate anxiety. The parameter  $Z$  captures the utility loss due to a shock and  $q$  is the likelihood of a shock conditional on seeing a shock pattern. The parameter  $r_z \in \{0, 1\}$  reflects whether shocks (hence, the subscript  $z$ ) are associated with right-tilted ( $r_z = 1$ ) or left-tilted ( $r_z = 0$ ) patterns in a given trial. In our model, anticipatory utility is linear in beliefs.<sup>14</sup>

The agent will not only experience the disutility of anticipatory anxiety, but also the disutility of actually receiving the shock, which is given by  $(r_z\hat{p} + (1 - r_z)(1 - \hat{p}))qZ$ .

Suppose next that self-deception is not frictionless, but instead subject to a quadratic cognitive cost  $\lambda(s)(p - \hat{p})^2$ . The cognitive cost function is increasing in the distance between a participant's initial belief and her chosen belief.  $\lambda$  captures the magnitude of the cognitive cost and we assume that  $\lambda$  is increasing in  $s$ , the strength of the signal the agent encounters. Then, the agent's total utility is given by

$$\begin{aligned} U = & \frac{1}{2} (1 + 2p\hat{p} - \hat{p}^2) M \\ & - (r_z p + (1 - r_z)(1 - p))qZ - \sigma_z(r_z\hat{p} + (1 - r_z)(1 - \hat{p}))qZ \\ & - \lambda(s)(p - \hat{p})^2 \end{aligned}$$

Maximizing the above expression with respect to  $\hat{p}$  yields a participant's optimal belief

$$\hat{p}^* = p(s, r_\theta) - \frac{\sigma_z(2r_z - 1)qZ}{M + 2\lambda(s)}$$

---

<sup>13</sup>The BDM mechanism was used in experiment 1 whereas experiments 2-4 paid participants for accurately identifying a given pattern. We cast the model in terms of the BDM mechanism for its analytical convenience. After Experiment 1 established similar results for confidence and (binary) accuracy judgments, we implemented discrete incentives in the subsequent online experiments in order to shorten instructions and reduce cognitive load.

<sup>14</sup>Some authors have assumed that anticipatory utility is concave in beliefs (e.g. Caplin and Leahy 2001), which implies that anxiety is then also induced by a greater variance in beliefs. The effect of non-linear belief-based utility on wishful thinking may be a fruitful topic of exploration for future work. What matters for our simple predictions here is that utility is monotonic in beliefs, we do not require linearity.

From this optimal belief we can derive hypotheses about the effects of our experimental treatments. We consider the case in which the true pattern is right-tilted,  $r_\theta = 1$ , so that  $\hat{p}$  is the belief in the correct answer. The case of  $r_\theta = 0$  is symmetric. Then, the *Shock* condition corresponds to  $r_z = 1$  and the *No-Shock* condition corresponds to  $r_z = 0$ . The amount of wishful thinking is given by

$$W := \hat{p}^*(r_z = 0) - \hat{p}^*(r_z = 1) = \frac{2\sigma_z q Z}{M + 2\lambda(s)} \quad (1)$$

From (1), and under the assumption that  $\sigma_z$  and  $\lambda$  are positive, we derive the following main hypothesis.

**Hypothesis 1 (Wishful thinking)** *There is positive wishful thinking, i.e.  $W > 0$ .*

Next, the effect of ambiguity on wishful thinking follows directly from our assumption that  $\lambda'(s) > 0$ .

**Hypothesis 2 (Ambiguity)** *Wishful thinking decreases when the pattern is easier to identify, i.e.  $\frac{dW}{ds} < 0$ .*

The test of hypothesis 2 illuminates how signal precision affects the production of distorted beliefs or a participant’s ability to self-deceive. Signal precision  $s$  also affects  $p(s, r_\theta)$ , which in turn impacts the motivation to hold distorted beliefs. But because, owing to our symmetric design,  $p(s, r_\theta)$  drops out of our measure of wishful thinking, we can shed the light on participants ability to self-deceive net of the strength of motives they may have to hold certain beliefs.

Next, the model predicts that higher accuracy incentives  $M$  raise the material costs of biased beliefs and make them less desirable.

**Hypothesis 3 (Incentives)** *Wishful thinking declines in the size of the accuracy bonus, i.e.  $\frac{dW}{dM} < 0$ .*

Experiment 2 varies psychological stakes by varying the loss associated with a loss pattern. By relabelling  $Z$  to capture this monetary loss, we can state the following hypothesis.

**Hypothesis 4 (Loss size)** *Wishful thinking increases in the disutility of the adverse outcome, i.e.  $\frac{dW}{dZ} > 0$ .*

Appendix C features a number of extensions of the model. First, in Section C.1, we show that the above predictions are robust to allowing the agent to derive anticipatory utility from her expectation of future accuracy payoffs. We show that the agent’s wishful thinking increases in this case, because her beliefs will be less disciplined by accuracy incentives. She now cares about the

actual *and* the perceived accuracy payoffs she obtains and the latter are not decreasing in biased beliefs.

Second, in Section C.2, we allow for a “bracing” or “defensive pessimism” motive for biased beliefs. We assume that, holding the actual likelihood of the shock constant, an agent suffers less disutility from the shock if she expects the shock to occur with a higher likelihood. We show that defensive pessimism works in the opposite direction of wishful thinking, so our main hypothesis can be rephrased as saying that wishful thinking trumps defensive pessimism as the dominant motive for belief distortion.

We can also use the model to account for correlations between measures of wishful thinking and (realized) anxiety, based on the interpersonal heterogeneity of fundamental parameters. In Section C.3 we show that this correlation can be positive, negative or zero depending on which underlying heterogeneity is driving differences in wishful thinking and anxiety. In particular, heterogeneity in  $\lambda$  implies a negative correlation and heterogeneity in  $\sigma_z$  implies a positive correlation. The empirical correlation therefore sheds light on the relative importance of different heterogeneities.

Our data confirms some predictions of the model and is at odds with some others. To capture these discrepancies, section C.4 proposes a revised model that allows for ex-ante investments in signal precision.

## 5 Results

We start with an overview of the main results across all of our experiments. The outcome of interest is whether pattern detection accuracy differs between shock/loss and no-shock/no-loss patterns. Table 2 shows the results of OLS models that regress accuracy on our treatment variables. To deal with interdependence between observations for a given participant, we take as a unit of observation the average accuracy over an individual’s trials within a given treatment and cluster standard errors at the participant level. The Appendix contains further data and analyses. In Appendix A we show that all of our results are robust to panel data regressions over all trials that include individual fixed effects and cluster standard errors at the individual level. Appendix Table A.1 provides descriptive statistics of accuracy levels for all of our experiments. Appendix Figure ?? provides an overview of the cumulative distribution functions of accuracy in shock/loss and no-shock/no-loss patterns.

Our main hypothesis is that participants are less accurate in identifying patterns associated with a shock or monetary loss. Table 2 exhibits strong evidence for this hypothesis. In all experiments, the coefficient for shock/loss patterns is negative and highly statistically significant. This also holds when we include interactions with the Difficulty and High Accuracy Bonus treatments.

The only qualification to this statement is needed in Experiment 1 (column 2), where we find statistical significance only for the more difficult patterns. Thus, wishful thinking appears as a robust phenomenon, both across sources of anxiety and across pattern recognition tasks.

We also hypothesize that wishful thinking is more pronounced for ambiguous or difficult patterns, where the signal is weaker and it may be easier to convince oneself of a positive outcome. The coefficient on the difficulty level across patterns shows participants are less likely to be correct on difficult patterns, where the size of the coefficient reflects how we operationalized difficulty in the different experiments (see Table 1 for details). Moreover, the interaction term shows that this is especially pronounced for loss or shock patterns, thus confirming our hypothesis in all experiments.

The third hypothesis is that incentives for accuracy reduce wishful thinking, as they raise the costs of wrong beliefs. Table 2 shows no evidence for this hypothesis, as the interaction terms between loss/shock pattern and the accuracy bonus are not statistically significant. However, a closer examination, in Section 5.3, reveals that accuracy incentives can have an effect in some settings. Finally, we find that loss size, which we varied in Experiment 2 has a positive effect on accuracy and increases wishful thinking, but in both cases the effects are not statistically significant.

In the sections below, we elaborate on the results of the individual experiments and develop additional insights and interpretations.

As Table 2 shows, effect sizes differ quite a lot across experiments. For instance, effect sizes are about four times larger in Experiment 2 than in Experiment 1. It is difficult to directly compare these effects, as there were several differences between these experiments. Besides replacing shocks with losses, Experiment 2 took place online, which necessitated changes to the exact instructions, the earning amounts and number of trials. Experiments 3 and 4 further differ in the perceptual task and other implementation details.

If we average wishful thinking over all participants in all experiments, then we find that average accuracy is 78.0 percent for no-loss/shock patterns and 69.7 percent for loss/shock patterns, respectively 28.0 percentage points and 19.7 percentage points above the 50 percent baseline of random choice.<sup>15</sup> Therefore, average wishful thinking is 8.3 percentage points and seeing a shock/loss rather than a no-shock/no-loss pattern decreases performance above chance by almost one third.

---

<sup>15</sup>We use as an observation the individual averages of accuracy for shock/loss and no-shock/loss patterns, so that every individual is weighted the same regardless of the number of trials in the experiment she completed.

Table 2: OLS regressions of accuracy levels on treatment across experiments

Dep var:	Experiment 1 (Electric Shocks)		Experiment 2 (Monetary losses)		Experiment 3 (Repeat flash)		Experiment 4 (Dot task)	
	(1) Accuracy	(2) Accuracy	(3) Accuracy	(4) Accuracy	(5) Accuracy	(6) Accuracy	(7) Accuracy	(8) Accuracy
Shock/Loss pattern	-4.111*** (1.264)	-2.014 (1.736)	-16.54*** (1.605)	-8.248** (3.489)	-4.266*** (0.766)	-3.052*** (0.865)	-8.453*** (1.040)	-7.393*** (1.308)
High accuracy bonus (HAB)	0.785 (0.878)	0.313 (1.387)	-0.588 (0.851)	-1.081 (1.089)	0.630 (0.474)	0.685 (0.601)	1.670*** (0.627)	1.004 (0.853)
Difficulty level (DL)	-8.602*** (0.634)	-7.318*** (0.795)	-15.68*** (1.019)	-11.04*** (1.114)	-20.55*** (0.668)	-19.39*** (0.794)	-7.073*** (0.269)	-6.497*** (0.361)
Loss size (LS)			-0.617 (0.906)	0.776 (1.245)				
Shock/Loss pattern x HAB		0.944 (1.787)		0.994 (1.771)		-0.110 (0.881)		1.330 (1.322)
Shock/Loss pattern x DL		-2.569** (1.102)		-9.200*** (1.701)		-2.317*** (0.892)		-1.150** (0.503)
Shock/Loss pattern x LS				-2.784 (1.869)				
Constant	80.75*** (1.106)	79.70*** (1.287)	85.82*** (1.964)	81.65*** (2.310)	87.66*** (0.791)	87.06*** (0.829)	89.53*** (0.732)	89.00*** (0.796)
Observations	720	720	3415	3415	3408	3408	6534	6534
$R^2$	0.261	0.266	0.134	0.140	0.236	0.236	0.110	0.110

OLS regressions of accuracy on treatment dummies and interactions. Each observation is the average accuracy of an individual over all trials in a given treatment. “Shock/Loss pattern” is a dummy if the pattern is associated with a shock (Experiment 1) or loss (Experiments 2-4). “High accuracy bonus” is a dummy that represents a high accuracy bonus, while “Difficulty level” is a categorical variable that counts the difficulty level of the perceptual task, with the number of levels dependent on the experiment (see Table 1 for details). The continuous difficulty levels in Experiment 3 were binarized using a median split. “Loss Size” refers to the size of the monetary loss that we varied in Experiment 2. Standard errors in parentheses clustered by individual. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

Our effect sizes are unlikely to be predictive of particular applications, since our results show a dependence on the task at hand and its difficulty. Nevertheless, as an external benchmark, one might consider mammogram reading, a complex pattern recognition task with a high-stake emotional outcome. Studies on interventions with radiologists often celebrate improvements in accuracy of a few percentage points, which are well in range of our effect sizes (Hadjiiski et al., 2004; Houssami et al., 2004).

## 5.1 Experiment 1: Electric Shocks

Figure 2 shows the average accuracy levels from Experiment 1, split by shock and no-shock patterns. Each observation is the individual average over all trials in a given category, so  $N = 60$  in each category.<sup>16</sup> Panel (a) compares average accuracy between shock and no-shock patterns, demonstrating wishful thinking of about four percentage points (72.3 vs. 68.6 percent). Panel (b) displays the impact of the three difficulty levels, as defined by the size of the tilt of the pattern. There appears to be some wishful thinking for easy patterns (2.4 percentage points) and medium patterns (2.5 percentage points). However, Table A.2 provides interaction terms for each of the difficulty levels, and shows that wishful thinking is statistically significant only for the most difficult patterns, where it rises to about 8 percentage points. Finally, panel (c) displays the impact of raising the prize for the BDM mechanism from 1 to 20 euro. Wishful thinking is about 1.4 percentage points more pronounced under the low bonus than under the high bonus, but the difference between the two conditions is not statistically significant.

**Confidence.** In addition to the accuracy measure, we elicited a measure of confidence in having correctly identified the pattern, incentivized with a BDM mechanism. Since we know the actual state of the world, this allows us to construct the variable “Belief”, which measures the subjective belief in the correct answer and provides a more continuous measure of participants’ perceptions. Beliefs vary on a scale from 0 (meaning the subject indicated 100% confidence in the wrong answer) to 100 (meaning the subject indicated 100% confidence in the correct answer). Figure A.1 and Table A.3 in Appendix A show results for this belief variable that are analogous to those for accuracy. We find the effects for accuracy and confidence are comparable both in size and in statistical significance.

---

<sup>16</sup>Table 2, in columns 1 and 2, provides the regression analysis associated with these results, and Appendix Table A.2 adds robustness across regression specifications. Moreover, in Appendix B, we provide a replication of Experiment 1 with  $N = 50$ .

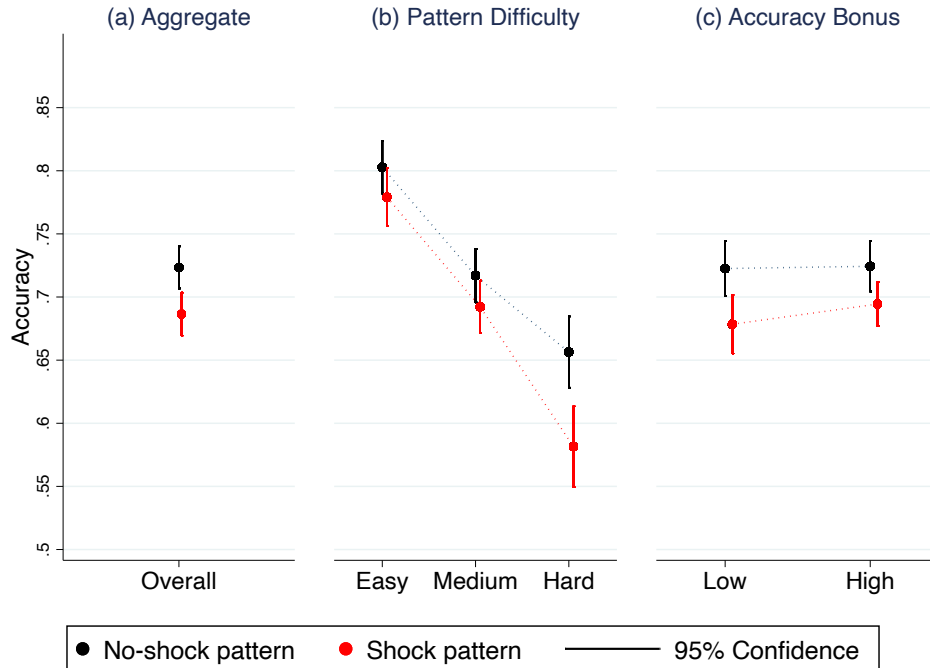


Figure 2: Electric shocks and accuracy in Experiment 1. Average accuracy levels, split by shock and no-shock pattern. Bars indicate 95% confidence intervals. One observation is the average over an individual’s trials in a given category, so  $N = 60$  in each category. Panel a) shows aggregate results. Panel b) disaggregates the results by difficulty (tilt) of the pattern. Panel c) disaggregates by incentives for accuracy.

## 5.2 Experiment 2: Monetary losses as a source of anxiety

We now turn to Experiment 2, which replaced electric shocks with monetary losses. While the literature has documented how the threat of electric shocks increases anxiety, no such evidence is available for losses. As a manipulation check, we therefore asked subjects to report their agreement with the statement “I felt anxious about losing money from my endowment” on a scale from 1 to 5 after each treatment block of 8 trials in which losses could occur. Panel (a) of Figure 3 shows the outcome of this manipulation check, where we count the scores in each block in both the Low Loss (10 cents) and High Loss (5 pounds) condition. This shows a clear difference between the two loss conditions, with average anxiety being 3.39 in the Low Loss condition and 4.15 in the High Loss condition ( $p < 0.001$  on a linear regression with standard errors clustered by participant). It also demonstrates that participants report substantial levels of anxiety about monetary losses even in the Low Loss condition.

Turning to the main results, Figure 4 shows the average accuracy levels from Experiment 2, split by Loss and No-loss patterns. Each observation is an individual’s average over all trials in a given category, so  $N = 221$  in each category. Table 2, columns 3 and 4, provides regression evidence



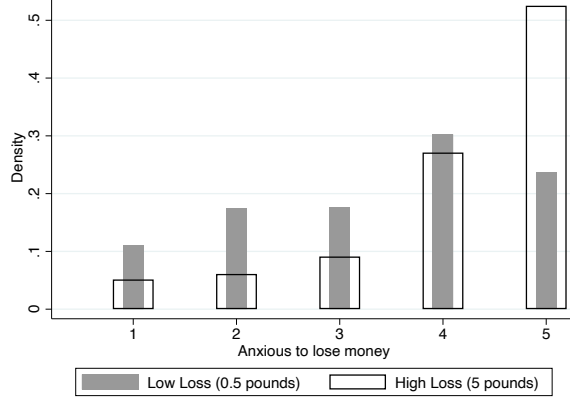


Figure 3: Manipulation check. Histogram of agreement with the statement “I felt anxious about losing money from my endowment” measured on a 5-point Likert scale, split by loss size. Each report in a treatment block counts as one observation.

associated with these results, and Table A.4 provides robustness across regression models. Results exclude the Neutral condition, since this is not a test of wishful thinking and is discussed in Section 5.5.3.

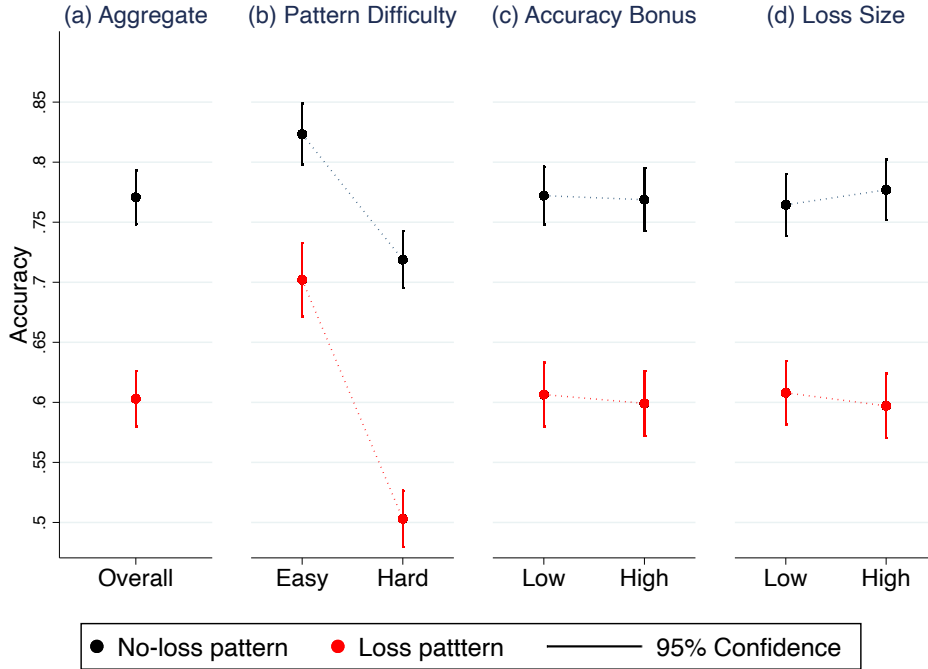


Figure 4: Monetary losses and accuracy in Experiment 2. Average accuracy levels, split by loss and no-loss patterns. Bars indicate 95% confidence intervals. One observation is the average over an individual’s trials in a given category, so  $N = 221$  in each category. Panel a) shows aggregate results. Panel b) disaggregates the results by difficulty (tilt) of the pattern. Panel c) disaggregates by incentives for accuracy. Panel d) disaggregates by size of losses.

Panel (a) of Figure 4 compares average accuracy on No-loss patterns with the Loss patterns. We see wishful thinking of 17 percentage points, which is highly statistically significant. The effect size is strikingly large: compared to the random-choice benchmark of 50 percent accuracy, the improvement in accuracy is almost 3 times higher under patterns associated with no loss compared to those that are associated with a loss. Panel (b) shows clear wishful thinking for both pattern difficulty levels, as well as an interaction effect between wishful thinking and difficulty. Panel (c) shows the effect of seeing a loss pattern for accuracy bonuses of 0.1 and 8 pounds respectively. We find a clear evidence for wishful thinking in both cases. However, there is no evidence that incentives improve performance, and no evidence for an effect of the accuracy bonus on wishful thinking, with wishful thinking being 16.6 percentage points under the low bonus and 17.0 percentage points under the high bonus. We will come back to the interpretation of this null result in Section 5.4. Finally, Panel (d) shows the effect of changing the loss size from 10 cents to 5 pounds. While this raises wishful thinking by about 2.7 percentage points, this difference is not statistically significant.

### 5.3 Experiments 3 and 4: Task characteristics

Figure 5 shows the average accuracy levels in Experiments 3 and 4, split into the Loss and No-loss conditions. As before, each observation is the individual average over all trials in a given category. Columns 4-8 of Table 2 provide regression evidence associated with these results and Tables A.5 and A.6 provide evidence for robustness across regression models.

**Results Experiment 3.** Subfigure (i) of Figure 5 shows the average accuracy levels from the sequential Gabor Task used in Experiment 3. Panel (a) compares average accuracy on no-loss patterns with the loss patterns. This results in wishful thinking of 4.4 percentage points. Task difficulty was a continuous variable in this task, defined by the posterior likelihood ratio of the two pattern-generating processes. Figure 5 Panel (b) displays the impact of a median split on this variable, and shows both a clear and statistically significant effect of higher difficulty on wishful thinking. Panel (c) shows wishful thinking for accuracy bonuses of 0.05 and 10 pounds respectively. Again, we find little evidence that incentives improve performance: A high bonus improves accuracy by about 0.7 ppt, but the effect is not close to being statistically significance. Moreover, there is no interaction with the loss pattern, so no reduction in wishful thinking from higher accuracy incentives.

**Results Experiment 4.** Subfigure (ii) of Figure 5 shows the average accuracy levels from the Dot-counting task used in Experiment 4. Panel (a) shows wishful thinking of 8.5 percentage points

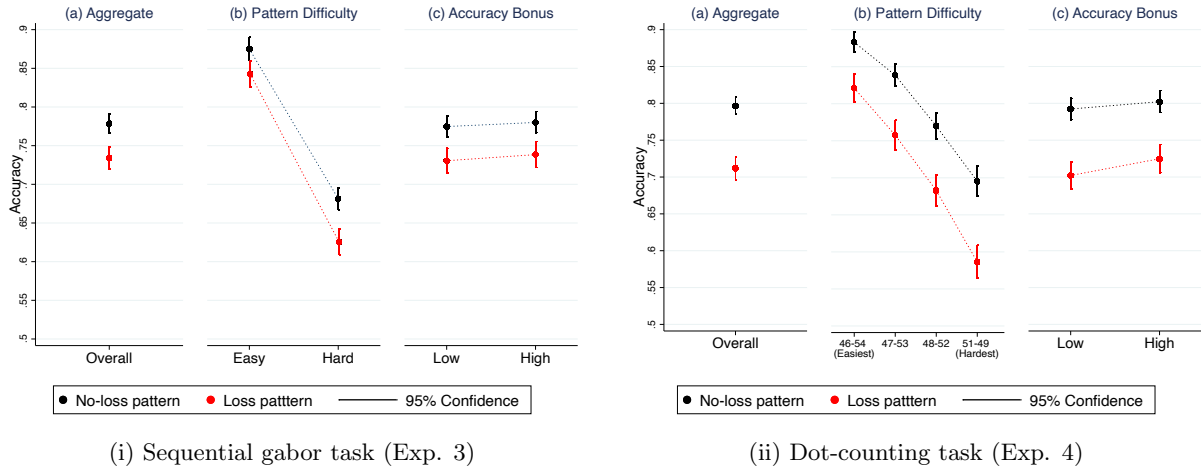


Figure 5: Accuracy in the multiple gabor and dot-counting tasks in Experiment 3 and 4. Average accuracy levels, split by loss and no-loss pattern. Bars indicate 95% confidence intervals. One observation is the average over an individual’s trials in a given category. In each subfigure, Panel a) shows aggregate results. Panel b) disaggregates the results by difficulty of the pattern, with a median split shown for Experiment 3. Panel c) disaggregates by incentives for accuracy.

in this task. Panel (b) displays the impact of pattern difficulty, where the easy patterns had a 46-54 split in colored dots, and the hardest patterns a 49-51 split. Once again, we confirm a statistically significant effect of difficulty on accuracy as well an interaction with wishful thinking. Panel (c) shows the pattern for the different levels of the accuracy bonus of 0.05 and 10 pounds. Unlike for the tasks we considered above, incentives improve performance: A high bonus improves accuracy by about 1.6 ppt, which is significant at the 5 percent level in the regression in Table A.1. However, as in our previous experiments, there is little interaction with the loss pattern. As we discuss next, this result hides important heterogeneities.

## 5.4 Incentives for accuracy

Across our experiments, we find no statistically significant effect of the accuracy bonus on wishful thinking. This may indicate that subjects do not make cognitive trade-offs between psychological and material incentives. Alternatively, it may be because participants did not care about or pay attention to the accuracy bonus or somehow were not able to react to it. We now discuss these explanations in more detail, and find that our aggregate results hide some important heterogeneity.

**Incentives and cognitive effort.** To investigate whether participants cared about, paid attention to, and responded to the bonus, we asked participants in Experiments 2, 3 and 4, at the

Table 3: Regressions of cognitive effort on accuracy bonus across experiments

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
	Concen- tration	Concen- tration	Concen- tration	Response time (log)	Response time (log)	Response time (log)	Response time (log)
High acc. bonus	0.117*** (0.0330)	0.120*** (0.0169)	0.172*** (0.0258)	0.0377** (0.0169)	0.0525*** (0.0129)	0.0311*** (0.00646)	0.130*** (0.0176)
Experiment no.	2	3	4	1	2	3	4
Fixed effects				✓	✓	✓	✓
Observations	442	852	818	11520	11396	33507	21217
$R^2$	0.007	0.017	0.012	0.001	0.003	0.001	0.013

Regressions of cognitive efforts on a dummy for the high accuracy bonus by experiment. Columns 1-3 show regressions on self-reported concentration, where an observation is an individual’s average concentration over all trials in the High Bonus and Low Bonus condition. Concentration is measured as agreement with the statement “In the past 8 trials I was very concentrated on the task” on 5-point Likert scale. Columns 4-7 show panel regressions where the outcome variable is log response time in each trial, where the latter is measured in milliseconds. Panel regressions include individual fixed effects. Standard errors in parentheses are clustered by individual. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

end of each 8 trial block, to report their agreement with the statement “I was very concentrated on the task”. Columns 1-3 of Table 3 show regressions that identify the increase in this variable resulting from a higher accuracy bonus, where each observation is an individual’s average reported concentration for each accuracy bonus level. In each of the experiments, a higher accuracy bonus leads to a significant increase in self-reported concentration.

In addition, if the high bonus leads to a more careful evaluation of answers, one would expect longer response times. Because of the highly the skewed nature of the response time distribution, which may be sensitive to large outliers, we look at the logarithm of response times as an outcome variable, which is measured in milliseconds. Columns 4-7 of Table 3 show panel regressions where we regress log response times on the high accuracy bonus, revealing an increase in response times. The accuracy bonus significantly increases response times in all experiments (we find similar results if we take raw response times), particularly in Experiment 4, which we will investigate in more detail below.

These results imply that participants care about and react to the accuracy bonus, albeit not by adjusting their wishful thinking. We can look at experiment 2, which featured the most wishful thinking of all experiments, to calculate an upper bound of the monetary cost associated with this stubborn wishful thinking. We zoom in on trials in which the loss and the correct answer were aligned, which mirror the many applications where what is true and what scares us aligns.<sup>17</sup>

<sup>17</sup>Ex-post, the symmetric nature of the task means that sometimes wishful thinking decreases accuracy (when losses are associated with the correct answer) and sometimes it increases accuracy (when losses are associated with

Comparing accuracy on such loss patterns in the High Bonus condition with accuracy in a set of control trials in which no losses were possible, implies an expected monetary cost from wishful thinking of about 87 cents.<sup>18</sup> This corresponds to roughly 10 minutes of work on the Prolific platform.

**Incentive effects and cognitive control.** The only time we find a statistically significant effect of accuracy incentives on accuracy is in the context of the dot counting task in Experiment 4. This experiment also sees the largest increases in our measures of cognitive effort in Table 3. This makes sense, as this task was chosen to be very elastic to cognitive effort. Participants who are really motivated to get the answer right may even *count* the dots. Incentive effects on performance and wishful thinking may therefore be particularly pronounced if they motivate participants to engage in such a strategy.

To explore this, we asked participants in the post-experimental questionnaire whether they counted dots. On this question 9% of subjects replied with “Always”, 38% replied “Sometimes”, and 53% replied “Never”. These answer categories do indeed correlate with participants’ response times. The participants in these three answer categories have mean response times of 14.4 seconds, 6.0 seconds, and 3.1 seconds respectively. Moreover, the participants in the Sometimes category also show a 32 percent increase in mean response times in reaction to the high accuracy bonus. This effect is far larger than that observed for participants in the Never (7 percent) and Always (14 percent) categories. Given the highly skewed nature of response time distributions appendix Table A.8 also analyzes the effect of the bonus on median response times. This shows that the increase is not significant for the Never category, but highly significant in the Sometimes category ( $p < 0.0001$ ). Moreover, it is marginally significant for the Always category and highly significant if we combine the Sometimes and Always categories.

Since the Sometimes counters are most responsive to the accuracy bonus, we investigate whether this category is also more likely to exhibit a reduction in wishful thinking. Figure 6 shows wishful thinking across the three counting categories. The figure shows clear evidence for a reduction in

---

the incorrect answer). As a result, averaged over all trials, the presence of losses does not decrease accuracy. This does not mean that wishful thinking is a money maximizing strategy from the subjective perspective of the agent. For an unbiased participant who is unsure which pattern she saw, self-deception always has negative expected value. This is true regardless of whether the bias pushes towards less accurate answers (for shock patterns) or more accurate answers (for no-shock patterns), because the agent’s only way to distinguish between these is her (initial) subjective belief.

<sup>18</sup>In the High Accuracy Bonus condition, participants could earn 8 pounds if their answer in a randomly selected trial belonging to that category was correct. In that condition, accuracy for loss patterns was 60.3 percent. Accuracy in trials that rule out any wishful thinking was 71.2. So across trials, associating the true state of the world with an anxiety-inducing outcome lead to a 10.9 ppt decrease in accuracy and an expected loss of  $0.109 * 8 = 0.87$  pounds. For the most ambiguous patterns, the decrease in accuracy is 12.8 ppt, and the expected loss is 1.02 pounds.

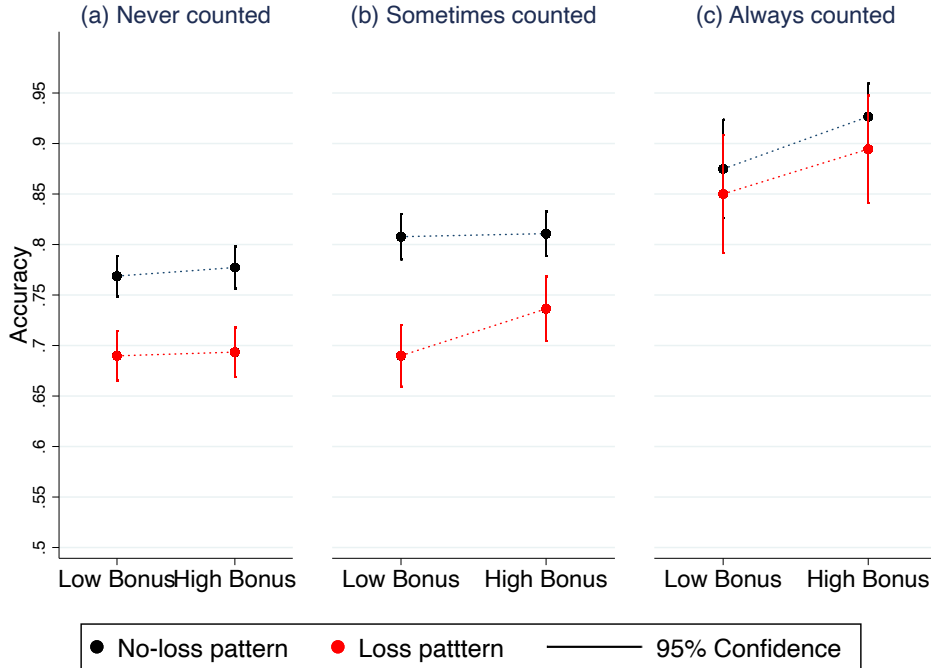


Figure 6: Accuracy in the dot-counting task. Average accuracy levels, split by loss and no-loss patterns. Bars indicate 95% confidence intervals. One observation is the average over an individual’s trials in a given category. Panel a) shows participants who report that they never count ( $N = 214$ ). Panel b) shows participants who sometimes counted ( $N = 154$ ). Panel c) shows participants who always counted ( $N = 35$ ).

wishful thinking among the Sometimes category. Appendix Table A.7 shows that this interaction is significant at the 5 percent level for that category, and marginally significant for all counters. The Never category show a slightly negative and insignificant interaction. The variation in accuracy between the categories of counters provides further evidence that effort reduces the scope for wishful thinking: the Always counters find the correct answer substantially more often than other participants and also show a highly reduced, and statistically not significant level of wishful thinking.

The results in this section provide evidence that incentives for accuracy can reduce wishful thinking. However, such an effect only obtains for participants who increase their accuracy through cognitive effort. The increase in accuracy may then act much like the exogenous decrease in pattern difficulty we saw affecting wishful thinking. By contrast, the self-deception we see taking place in all of our experiments does not seem to respond to material incentives directly, i.e. by participants calibrating the intensity of their wishful thinking or ex-post signal distortion to trade-off psychological and material incentives at the margin.

How do these results relate to our model in Section 4? The model predicts that a higher accuracy bonus reduces wishful thinking, by affecting the optimization process behind motivated

beliefs. The prediction is consistent with the behavior of counters in Experiment 4. Thus, in settings that allow for successful investments in signal precision, the model predicts correctly, even if only in an “as-if” sense. However, the model predictions fail in the setting of Experiments 1, 2 and 3, which precluded the possibility of improving accuracy, suggesting that accuracy incentives do not affect the ex-post signal distortion featured in the model.

In Appendix section C.4, we revise the model to be more in line with our findings. There, we assume that self-deceptive efforts are costless up to a certain point, but impossible thereafter. One interpretation is that self-deception is closer to an “automatic” or “system 1” process (see also Kappes and Sharot, 2019). This model implies that wishful thinking is slow to respond to psychological and material incentives at the margin. Successful investments in signal precision, on the other hand, can constrain wishful thinking by lowering the maximum possible amount of self-deception. Higher accuracy bonuses increase such investments, irrespective of whether agents are sophisticated about the effect of signal precision on wishful thinking, and, in line with our data, then reduce wishful thinking whenever they are successful.

## 5.5 Robustness and competing explanations

In this section we conduct a number of robustness checks and investigate competing explanations for our results.

### 5.5.1 Robustness

At the end of the experiment, we asked subjects several questions about their perceptions of the experiment. We use these variables to conduct several robustness checks that could identify the potential effects of confusion, misunderstanding, or distrust in the experimenter on wishful thinking. In these analyses, we pool the data from all experiments. In a first check, we restrict our sample by excluding participants who scored high on perceived difficulty of the instructions, a general measure of understanding. In further robustness checks we exclude those who found it hard to recall the treatment conditions, who made more than two mistakes in the initial control questions, who did not trust the experimenters, or those whose accuracy in the experimental task was below 60 percent. The latter measure excludes some participants who answer almost randomly and a small number of participants who almost always select the no-shock pattern.

The results are reported in Appendix Table A.9: wishful thinking remains highly significant in all selected samples, with small and statistically insignificant changes in effect sizes. The interaction of shock patterns with pattern difficulty also remains statistically significant in all specifications. The estimate for the interaction effect between the accuracy bonus and the shock pattern is generally

positive but not statistically significant. Table A.10 shows similar results in analogous regressions where we use panel data from all trials and include individual fixed effects. We conclude that our results are not driven by misunderstanding or distrust.

### 5.5.2 Illusion of control

Our experimental instructions stress that participants' answers do not have a causal effect on the shocks or losses. Several quiz questions during the instruction phase explicitly asked subjects to confirm their understanding of this point. Nevertheless, participants may have somehow come to believe during the experiment that their answers were associated with shocks or losses. Such an 'illusion of control' may lead subjects to switch their answers to the non-aligned pattern in order to avoid the shock or loss.

To address this point, we conducted another understanding check in the closing questionnaire of Experiments 2, 3 and 4. A multiple choice question asked participants what drove losses in the experiment: a) the tilt of the pattern and designated loss category, b) their own answers, c) both, or d) don't know. On this question, the 77 percent of subjects who correctly gave the first answer had an average wishful thinking of 8.32 percentage points, while those who selected one of the other answers had indistinguishable average wishful thinking of 8.30 percentage points ( $p = 0.99$ , t-test). In column 7 of Appendix Tables A.9 and A.10, we also run our main regressions without the participants who answered the control question incorrectly. We find that the estimated effect size for wishful thinking is statistically and quantitatively robust.

### 5.5.3 Does seeing a shock or loss pattern increase noise?

It is possible that seeing a pattern that is associated with a loss or shock increases noise in their answer, thereby reducing accuracy for shock patterns. This "noise-based explanation" supposes that participants perceive the correct answer initially, but that the anxiety from observing a shock pattern reduces performance through some form of interference that differs from wishful thinking.

However, this alternative account cannot explain our data. First, the noise-based explanation would predict higher effect of shock/loss threat for easier patterns, because these induce a higher subjective probability of seeing a shock pattern and, hence, higher noise. However, we see the reverse in the data. Second, the noise-based explanation predicts that average accuracy should increase in a neutral condition where there is no threat of a shock or loss at all. Performance in such an anxiety-free condition should exceed those on shock patterns as well as the aggregate performance under shock and no-shock patterns. Note that it need not be higher than the performance under no-shock patterns, as in this case self-deception goes in the direction of the correct answer and



increases accuracy relative to neutral patterns.

To test this prediction, we conducted a Neutral treatment in both Experiment 2 and Experiment 4. In one part of the experiment, implemented in random order, subjects were informed that they could not lose money from their endowment in any trial of this part. We compare accuracy for neutral patterns with accuracy for loss and no-loss patterns, where we pool the data from the two loss sizes in Experiment 2. As before, we take as an observation the individual accuracy rate in each of these conditions. In both Experiment 2 and 4, we find that average accuracy for neutral patterns is between that of the loss and no-loss patterns. In Experiment 2, respective accuracy rates are 71.2 percent (neutral), 60.3 percent (loss), and 77.0 percent (no-loss). In Experiment 4, the corresponding percentages are 75.7 (neutral), 71.1 (loss) and 79.7 (no-loss). All within-experiment differences are statistically significant in a regression analysis (see appendix Table A.11). Furthermore, there is not much evidence that stress impedes average performance: accuracy is slightly (2.7 percentage points) higher in the Neutral condition than the average of the Loss and No-loss condition in Experiment 2, but not in Experiment 4, where they are almost identical (see Appendix Table A.1). Finally, a Neutral treatment in the replication of Experiment 1 further confirm these patterns, details of which are in Appendix B. We conclude that the data reject the noise-based explanation.

#### 5.5.4 Dynamics

Our experiments consist of many trials and within-subject treatments, so we can ask how wishful thinking evolves over time. It may be the case that participants get desensitized to the anxiety-inducing effects of electrical and monetary shocks as they experience an ever greater number of trials. They may then exhibit less wishful thinking in later trials. It may also be the case that initial experiences with losses or shocks heighten subsequent anxiety and increase wishful thinking in later trials. Our dataset offers a window into how motivated beliefs respond to experience and speaks to mechanisms that may be at play in real world settings, which often feature dynamics and an element of repetition.

In Appendix Figure A.5 we provide a visual overview of wishful thinking over time in each experiment. Appendix Table A.12 analyses statistically how the effect of seeing a loss or shock pattern on accuracy (our measure of wishful thinking) evolves over time by interacting a dummy for whether a participant sees a loss pattern with the number of trials a participant has gone through. In a second set of analyses we simply compare wishful thinking in the first half and the second half of the experiment. In Experiment 1, wishful thinking in the first half of the experiment is more than twice as large as wishful thinking in the second half. The coefficient is just shy of

statistical significance at conventional levels ( $p = 0.102$ ), but suggestive of the idea that participants get distracted or desensitized to the anxiety-inducing effects of electric shocks as trials go by. On the other hand, in Experiment 3, which features monetary losses, wishful thinking is higher in later trials and in the second half of the experiment. There is no significant effect of time or experience on wishful thinking in Experiments 2 and 4.

We can also ask how the experience of a shock or loss in the previous trial affects wishful thinking in the current trial. Appendix table A.13 investigates the effect of lagged shocks or losses on wishful thinking. We see no effect in any of the four experiments, regardless of whether or not we control for the time trend in wishful thinking.

While the structure of our dataset allows us to analyse the evolution of wishful thinking with experience, no unambiguous story emerges. The data suggest that repeated shocks may lead to desensitisation, but that experience with monetary losses can lead to heightened wishful thinking in some contexts. This suggests that the source of anxiety may matter for the dynamics of wishful thinking.

## 6 Wishful thinking as a trait

Motivated cognition, of which wishful thinking is one example, is usually identified by inducing experimental variation in participants' motives to hold biased beliefs. Since this experimental variation tends to be administered between subjects, the literature has not been able to obtain individual measures of a proclivity for motivated cognition and, with a few exceptions noted in Section 2, has therefore not been able to say much about individual differences. Instead, our within-subject design with many trials allows us to compute individual measures of wishful thinking and ask two underexplored questions: are there individual differences in wishful thinking and, if so, does wishful thinking correlate with other individual characteristics?

### 6.1 Wishful thinking as a personal trait

Do individuals differ in their proclivity for wishful thinking? Appendix Figure A.3 depicts histograms of individual-level wishful thinking in each experiment. We see that there is substantial variance, with a majority of participants engaging in some wishful thinking and some participants exhibiting the opposite treatment effect.

To establish that this apparent heterogeneity is not merely driven by measurement error or other sources of noise, we test for the stability of wishful thinking within individuals. In particular, we ask whether a participant's wishful thinking measured in one half of trials correlates with their

Table 4: Half-split correlations

X/Y	Wishful thinking			Accuracy		
	(1) Exp. 2	(2) Exp. 3	(3) Exp. 4	(4) Exp. 2	(5) Exp. 3	(6) Exp. 4
Odd/even trials	0.641	0.461	0.570	0.592	0.730	0.476
Difficult/easy patterns	0.573	0.526	0.457	0.575	0.497	0.563
First/second half	0.441	0.435	0.350	0.663	0.568	0.478
Low/high losses	0.460	-	-	0.589	-	-

Correlations between individual participants’ wishful thinking or accuracy as measured in X and Y trials. X and Y correspond to odd and even, difficult and easy, first and second half, and low and high loss trials respectively.

wishful thinking in the other half. For this exercise we split trials into odd and even numbered trials, trials with easy and trials with difficult patterns, trials in the first half and trials in the second half of the experiment and, for Experiment 2, trials with high stakes and trials with low stakes. Calculating such half-split correlations is common in psychology, where they are used to assess the reliability of individual measures derived from cognitive tasks (for example, see Pronk et al. 2021).<sup>19</sup>

Columns 1 through 3 of Table 4 report the half-split correlations of Experiments 2, 3 and 4 respectively.<sup>20</sup> Correlations are around 0.5, with some fluctuations depending on how we split the data, indicating that heterogeneity in wishful thinking reflects individual differences. Moreover, our measure of wishful thinking is only slightly less reliable or stable than participants’ skill in the pattern recognition tasks, as shown by the half split correlations in accuracy that we report in columns 4 through 6 of Table 4. To further show that our results are not driven by a few outliers, Appendix Figure A.4 shows the scatterplots pertaining to the odd-even trial splits in Table 4.

## 6.2 Emotional and cognitive covariates of wishful thinking

Since wishful thinking appears to be a stable individual characteristic, a natural next question to ask is whether it correlates with other variables, and whether such correlations can help understand the drivers of wishful thinking. First, we look at a self-reported measure of concentration, which we measured in the interblock surveys of Experiment 2, 3 and 4. Increased concentration may lead to more precise perceptions and higher accuracy, which in turn constrains wishful thinking. This

<sup>19</sup>In the same vein, our results here also help us assess how reliably our experimental design identifies wishful thinking.

<sup>20</sup>We exclude Experiment 1 because there we recalibrated both the strength of the shock and the difficulty of the patterns during the experiment. This confounds the half-split correlations of wishful thinking and accuracy.

is the mechanism suggested by our results on dot counters in Experiment 4.

A second covariate of interest is “defensive pessimism”, which measures the degree to which people adopt pessimistic beliefs to avoid disappointment.<sup>21</sup> This belief-based utility motive for self-deception into more pessimistic beliefs may arise if people are loss averse over changes in beliefs, as in Kőszegi and Rabin (2009).<sup>22</sup> Defensive pessimism runs counter to wishful thinking, as we show formally in Appendix C.2, so one would expect a negative correlation.

Finally, we are interested in the correlation of wishful thinking with self-reported anxiety about losing money from the endowment, which we measured in the interblock survey in Experiments 2 and 4 and hypothesized as a key antecedent of wishful thinking. The sign of this correlation is theoretically ambiguous, as we explain formally in Appendix Section C.3. If participants vary strongly in their ability to self-deceive, then higher wishful thinking should be associated with lower experienced anxiety. Conversely, if the primary source of heterogeneity between participants’ wishful thinking is their proneness to anxiety, then wishful thinking should be positively correlated with experienced anxiety.

Table 5 shows OLS regressions of wishful thinking on these explanatory variables. To generate maximal statistical power, we pool the data from all experiments in which the relevant explanatory variables were elicited. All regressions contain experiment dummies to control for differences in wishful thinking that are based solely on differences in the experimental task. Column 1 shows that wishful thinking is negatively correlated with the average self-reported concentration on pattern recognition. Wishful thinking therefore appears to be constrained by cognitive effort, presumably because participants who concentrate more are able to generate significantly more accurate representations of the signal.<sup>23</sup>

The correlation between wishful thinking and defensive pessimism is negative and not statistically significant at conventional levels. In column 2 we add a participant’s average self-reported anxiety to the regression model. The regression excludes Experiments 1 and 3, where we did not elicit an anxiety report. Anxiety is positively correlated with wishful thinking, with significance at

---

<sup>21</sup>Our measure is based on the defensive pessimism questionnaire (Norem, 2008). Following Lim (2009), we focus on the pessimism sub-scale, which measures agreement with the following statements: 1. I often start out expecting the worst, even though I will probably do OK. 2. I worry about how things will turn out. 3. I often worry that I won’t be able to carry through my intentions. 4. I spend lots of time imagining what could go wrong. 5. I imagine how I would feel if things went badly. 6. In these situations, sometimes I worry more about looking like a fool than doing really well.

<sup>22</sup>The presence of both anticipatory utility motives and a desire to avoid disappointment also has implications for the likely time-path of beliefs (Macera, 2014). Unfortunately, we only observe static beliefs.

<sup>23</sup>Regressing participants’ accuracy on their self-reported concentration, while controlling for experiment fixed effects, yields a highly statistically significant coefficient ( $p < 0.001$ ): going up one step on the five point Likert scale on which we measured concentration increases accuracy by 2.5 percentage points.

Table 5: Emotional and cognitive covariates of wishful thinking.

	(1)	(2)	(3)	(4)
Dep. variable:	Wishful Thinking	Wishful Thinking	Wishful Thinking	Wishful Thinking
Concentration	-0.0289*** (0.00958)	-0.0345*** (0.0122)	-0.0397*** (0.0107)	-0.0498*** (0.0134)
Defensive pessimism	-0.00604 (0.00399)	-0.0107* (0.00609)	-0.00905** (0.00425)	-0.0149** (0.00678)
Anxiety		0.0157* (0.00826)		0.0198** (0.00891)
Constant	0.328*** (0.0494)	0.317*** (0.0640)	0.385*** (0.0559)	0.386*** (0.0733)
Experiment dummies	✓	✓	✓	✓
Restrictions	None	None	Difficult instructions <4 of 7	Difficult instructions <4 of 7
Observations	1049	624	743	421
$R^2$	0.066	0.054	0.086	0.077

OLS regressions of wishful thinking on cognitive covariates. Data are from experiments 2, 3 and 4 in columns 1 and 3 and from experiments 2 and 4 in columns 2 and 4. Columns 3 and 4 only include participants with one of the three lowest scores on the question “How difficult did you find it to follow the instructions of this experiment?” measured on a 7-point Likert scale from very easy to very difficult. All regressions contain experiment dummies. Standard errors in parentheses. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

marginal levels.<sup>24</sup> Furthermore, the negative correlation between defensive pessimism and wishful thinking is now also marginally significant. In columns 3 and 4 we do a robustness check on these correlations, akin to that in Section 5.5.1, and exclude participants who reported that they found following the instructions hard to follow. Excluding such potentially noisy participants results in stronger correlations between wishful thinking and all covariates, adding confidence to the results.

These results allow us to sharpen our interpretations of wishful thinking. First, the negative correlation with concentration further underscores how cognitive effort can constrain wishful thinking through its effect on accuracy. Second, the negative correlation with defensive pessimism suggests that belief-based utility motives that run counter to wishful thinking exist and can be detected in the cross-participant heterogeneity of belief biases. Since defensive pessimism is a self-reported survey scale, its correlation with wishful thinking suggests that people are at least somewhat conscious of their tendencies for probability distortion. Finally, the positive correlation with self-reported

<sup>24</sup>We also elicited Beck Anxiety Inventory (BAI), a more general measure of anxiety that screens for, among other things, frequent physical symptoms of anxiety. BAI correlates with our measure of self-reported anxiety about incurring monetary losses in the experiment (corr=0.3,  $p < 0.001$ ), thereby validating our more focused and tailor-made measure. However, perhaps not surprisingly, the positive correlation between BAI and wishful thinking is not statistically significant.

Table 6: Real world attitudes and beliefs

Dep. variable:	(1) Risk seeking	(2) After- life	(3) Climate worry	(4) Risk seeking	(5) After- life	(6) Climate worry
Wishful Thinking	0.489* (0.254)	0.365 (0.354)	0.0295 (0.267)	0.762** (0.319)	1.018** (0.448)	0.110 (0.348)
Constant	3.584*** (0.121)	3.416*** (0.169)	5.704*** (0.128)	3.577*** (0.159)	3.228*** (0.223)	5.675*** (0.174)
Experiment dummies	✓	✓	✓	✓	✓	✓
Restrictions	None	None	None	Difficult instr. <4 of 7	Difficult instr. <4 of 7	Difficult instr. <4 of 7
Observations	1007	1007	1007	724	724	724
$R^2$	0.004	0.003	0.018	0.009	0.010	0.025

OLS regressions of wishful thinking on real world beliefs and attitudes. Risk seeking was measured as the answer to the question "Are you rather a risk-taking or risk-averse person (trying to avoid risks)?" on a Likert scale from 1 (very risk-averse) to 7 (very risk-seeking). Afterlife and climate worry were measured as agreement with the following statements on a 7-point Likert scale: "I believe in the existence of an afterlife.", "I am worried about climate change". Data is from experiments 2, 3 and 4 in columns 1, 2, 4 and 5. Data is from experiments 2 and 4 in columns 3 and 6. Columns 3 and 4 only include participants with one of the three lowest scores on the question "How difficult did you find it to follow the instructions of this experiment?" measured on a 7-point Likert scale from very easy to very difficult. Standard errors in parentheses. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

anxiety suggest that people differ in their innate anxiety, and that these differences are not (fully) overcome by their wishful thinking.

### 6.3 Does wishful thinking explain real world attitudes and beliefs?

To see whether wishful thinking helps explain other attitudes and beliefs, Table 6 features OLS regressions of three unincentivized and self-reported outcomes: risk seeking, belief in an afterlife, and worry about climate change. Columns 1 to 3 include data from all participants from whom the relevant dependent variables were elicited. Columns 4 to 6 again restrict the sample to participants who found the instructions easy to understand.

In column 1, we find a positive and marginally significant correlation between wishful thinking and risk seeking and in column 2, a non-significant positive correlation between wishful thinking and belief in an afterlife. These correlations become more significant in the restricted sample of columns 4 and 5, where we exclude those who report having difficulty with the experimental instructions. In columns 3 and 6, worry about climate risks is not significantly correlated with wishful thinking. We note that our results on the marginally significant positive correlations are exploratory, subject to concerns about multiple hypothesis testing, and that they should be confirmed by future research.

## 7 Conclusion

Philosophers and economists have long considered the importance of beliefs for people’s well-being. Jevons (1879) argues that “the greatest force of feeling and motive arises from the anticipation of a long-continued future”, while Bentham (1789) points to expectation as being among the most significant sources of pleasure and pain. Over the last decades, economists have introduced anticipatory feelings as a source of utility into their formal models (Loewenstein, 1987; Caplin and Leahy, 2001) and the notion of utility from anticipation has experienced somewhat of a “renaissance” (Loewenstein and Molnar, 2018; Molnar and Loewenstein, 2021).

Our experiments show the importance of such anticipatory emotions for belief formation. In each of the four experiments, participants are significantly less accurate in identifying patterns that may result in adverse outcomes. Such wishful thinking is most pronounced when evidence is ambiguous, a result that replicates across tasks with distinct sources of ambiguity. We find evidence that a higher material cost of wrong beliefs can reduce wishful thinking, but only when accuracy in the inference task is elastic to effort, so that participants can obtain more precise representations of signals if they choose to. Whether motivated beliefs respond to material incentives more generally is therefore likely to depend on the inference task and context in which beliefs are formed. Finally, we show that individuals differ in their propensity to engage in wishful thinking, with some showing the opposite tendency that reflects defensive pessimism.

Our results speak to decision making in a wide range of applications, as anticipatory anxiety has been invoked in decisions related to health, insurance, finance and politics.<sup>25</sup> They help explain why people seek solace and comfort in religious beliefs, why financial professionals ignore red flags about their asset portfolio, why people most at risk of a disease sometimes avoid testing for it, and why voters who are concerned about their jobs and the future of their children are susceptible to false but reassuring political narratives. The crucial role of ambiguity gives a rationale for the avoidance of precise information such as that provided in medical tests and helps explain the persistence of beliefs in phenomena such as the afterlife that, by their nature, do not admit clear evidence. The subtle findings on the role of accuracy incentives indicate that the bias can persist despite high personal costs.

Our results point to number of open questions. For instance, it would be interesting to explicitly

---

<sup>25</sup>Examples are beliefs about health risks (Schwardmann, 2019), financial decisions and excessive trading (Brunermeier and Parker, 2005; Eisenbach and Schmalz, 2015; Bridet and Schwardmann, 2017), time inconsistency (Caplin and Leahy, 2001; Köszegi, 2010), occupational choice and the labor market equilibrium (Akerlof and Dickens, 1982; Santos-Pinto et al., 2018), information acquisition (Yariv, 2002; Eliaz and Spiegler, 2006; Loewenstein, 2006), principal-agent communication (Köszegi, 2006; Caplin and Leahy, 2004), self-image and taboos (Bénabou and Tirole, 2011), groupthink (Bénabou, 2013) and politics (Bénabou, 2008; Levy, 2014; Le Yaouanq, 2016).

compare wishful thinking in the loss and gain domains, as loss aversion suggest that motives for wishful thinking will be stronger in the former domain. More generally, it would be important to understand the conditions under which wishful thinking responds to differences in losses or psychological stakes at the margin. It would also be interesting to look at whether the option to take a (costly) action to avert the adverse outcome, i.e. empowerment, can reduce wishful thinking. Finally, our results on how wishful thinking responds to material incentives suggest that cognitive investments in the accurate representation of signals are a key mechanism in motivated belief formation. Theoretical models, which have by and large focused on how material and psychological incentives affect ex-post belief distortions, may benefit from taking this mechanism seriously.

## References

- Akerlof, George A. and William T. Dickens**, “The economic consequences of cognitive dissonance,” *The American Economic Review*, 1982, 72 (3), 307–319.
- Armor, David A and Aaron M Sackett**, “Accuracy, error, and bias in predictions for real versus hypothetical events.,” *Journal of personality and social psychology*, 2006, 91 (4), 583.
- Auriol, Emmanuelle, Julie Lassebie, Amma Panin, Eva Raiber, and Paul Seabright**, “God insures those who pay? Formal insurance and religious offerings in Ghana,” *Mimeo, Toulouse School of Economics*, 2017.
- Balcetis, Emily and David Dunning**, “See what you want to see: Motivational influences on visual perception,” *Journal of Personality and Social Psychology*, 2006, 91 (4), 612–625.
- Barron, Kai**, “Belief updating: Does the ‘good-news, bad-news’ asymmetry extend to purely financial domains?,” *WZB Discussion Paper*, 2016, (309).
- , “Belief updating: does the ‘good-news, bad-news’ asymmetry extend to purely financial domains?,” *Experimental Economics*, 2021, 24 (1), 31–58.
- Bénabou, R and J Tirole**, “Self-confidence and personal motivation,” *The Quarterly Journal of Economics*, 2002, (August), 871–915.
- Bénabou, Roland**, “Ideology,” *Journal of the European Economic Association*, 2008, 6 (May), 321–352.
- , “Groupthink: Collective Delusions in Organizations and Markets,” *The Review of Economic Studies*, 2013, 80 (2), 429–462.



- **and Jean Tirole**, “Identity, Morals, and Taboos: Beliefs as Assets,” *The Quarterly Journal of Economics*, 2011, *126* (2), 805–855.
- Bentham, Jeremy**, *An Introduction to the Principles of Morals and Legislation*, Oxford: Clarendon Press, 1789.
- Bentzen, Jeanet Sinding**, “In crisis, we pray: Religiosity and the COVID-19 pandemic,” *Journal of Economic Behavior & Organization*, 2021, *192*, 541–583.
- Berns, Gregory S, Jonathan Chappelow, Milos Cekic, Caroline F Zink, Giuseppe Pagnoni, and Megan E Martin-skurski**, “Neurobiological Substrates of Dread,” *Science*, 2006, *754* (May), 754–758.
- Bosch-Rosa, Ciril, Daniel Gietl, and Frank Heinemann**, “Risk-Taking under Limited Liability: Quantifying the Role of Motivated Beliefs,” *Available at SSRN 3985775*, 2021.
- Bracha, Anat and Donald J. Brown**, “Affective decision making: A theory of optimism bias,” *Games and Economic Behavior*, 2012, *75* (1), 67–80.
- Bridet, Luc and Peter Schwardmann**, “Selling Dreams: Endogenous Optimism in Lending Markets,” *Mimeo, Ludwig Maximilians Universität München*, 2017.
- Brunnermeier, Markus K. and Jonathan A. Parker**, “Optimal expectations,” *American Economic Review*, 2005, *95* (4), 1092–1118.
- Burton, Jason W, Adam JL Harris, Punit Shah, and Ulrike Hahn**, “Optimism where there is none: Asymmetric belief updating observed with valence-neutral life events,” *Cognition*, 2022, *218*, 104939.
- Buser, Thomas, Leonie Gerhards, and Joël J. Van der Weele**, “Responsiveness to Feedback as a Personal Trait,” *Journal of Risk and Uncertainty*, 2018, *56*, 165–92.
- Camerer, Colin F. and Robin M. Hogarth**, “The effects of financial incentives in experiments: A review and capital-labor-production framework,” *Journal of Risk and Uncertainty*, 1999, *19* (1), 7–42.
- Caplin, A and J Leahy**, “Psychological expected utility theory and anticipatory feelings,” *The Quarterly Journal of Economics*, 2001, (February 2001), 55–79.
- Caplin, Andrew and John Leahy**, “The supply of information by a concerned expert,” *Economic Journal*, 2004, *114* (497), 487–505.

- **and** –, “Wishful Thinking,” *NBER working paper 25707*, 2019.
- **and Mark Dean**, “Revealed preference, rational inattention, and costly information acquisition,” *NBER Working Papers, (19876)*, 2014.
- Coutts, Alexander**, “Testing Models of Belief Bias: An Experiment,” *Games and Economic Behavior*, 2019, *113*, 549–565.
- Dean, Mark and Nate Leigh Neligh**, “Experimental tests of rational inattention,” *Discussion paper Columbia University*, 2019.
- Dewan, Ambuj and Nathaniel Neligh**, “Estimating information cost functions in models of rational inattention,” *Journal of Economic Theory*, 2020, *187*, 105011.
- Drobner, Christoph**, “Motivated beliefs and anticipation of uncertainty resolution,” *American Economic Review: Insights*, 2022, *4* (1), 89–105.
- Drugowitsch, Jan, Valentin Wyart, Anne-Dominique Devauchelle, and Etienne Koechlin**, “Computational precision of mental inference as critical source of human choice suboptimality,” *Neuron*, 2016, *92* (6), 1398–1411.
- Dunning, David and Emily Balcetis**, “Wishful Seeing : How Preferences Shape Visual Perception,” *Journal of Experimental Social Psychology*, 2013, *22* (1), 33–37.
- Eisenbach, Thomas M. and Martin C. Schmalz**, “Anxiety, overconfidence and excessive risk taking,” *Staff Report no. 711, Federal Reserve Bank of New York*, 2015.
- Eliasz, Kfir and Ran Spiegler**, “Can anticipatory feelings explain anomalous choices of information sources?,” *Games and Economic Behavior*, 2006, *56* (1), 87–104.
- Engelmann, J. B., F. Meyer, E. Fehr, and C. C. Ruff**, “Anticipatory Anxiety Disrupts Neural Valuation during Risky Choice,” *Journal of Neuroscience*, 2015, *35* (7), 3085–3099.
- Engelmann, Jan B., Friederike Meyer, Christian C. Ruff, and Ernst Fehr**, “The neural circuitry of affect-induced distortions of trust,” *Science Advances*, 2019, *5* (3), eaau3413.
- Enke, Benjamin, Uri Gneezy, Brian Hall, David Martin, Vadim Nelidov, Theo Offerman, and Jeroen van de Ven**, “Cognitive Biases: Mistakes or Missing Stakes?,” *The Review of Economics and Statistics*, 2021, pp. 1–45.

- Eyal, Peer, Rothschild David, Gordon Andrew, Evernden Zak, and Damer Ekaterina**, “Data quality of platforms and panels for online behavioral research,” *Behavior Research Methods*, 2021, pp. 1–20.
- Falk, Armin and Florian Zimmermann**, “Beliefs and utility: Experimental evidence on preferences for information,” 2016.
- Findling, Charles and Valentin Wyart**, “Computation noise in human learning and decision-making: Origin, impact, function,” *Current Opinion in Behavioral Sciences*, 2021, 38, 124–132.
- Ganguly, Ananda and Joshua Tasoff**, “Fantasy and Dread: The Demand for Information and the Consumption Utility of the Future,” *Management Science*, 2016, (September 2017), mns.2016.2550.
- Garbers, Yvonne and Udo Konradt**, “The effect of financial incentives on performance: A quantitative review of individual and team-based financial incentives,” *Journal of occupational and organizational psychology*, 2014, 87 (1), 102–137.
- Grillon, Christian**, “Models and mechanisms of anxiety: evidence from startle studies,” *Psychopharmacology*, 2008, 199 (3), 421–437.
- Hadjiiski, Lubomir, Heang-Ping Chan, Berkman Sahiner, Mark A Helvie, Marilyn A Roubidoux, Caroline Blane, Chintana Paramagul, Nicholas Petrick, Janet Bailey, Katherine Klein et al.**, “Improvement in radiologists’ characterization of malignant and benign breast masses on serial mammograms with computer-aided diagnosis: an ROC study,” *Radiology*, 2004, 233 (1), 255–265.
- Haisley, Emily C and Roberto A Weber**, “Self-serving interpretations of ambiguity in other-regarding behavior,” *Games and economic behavior*, 2010, 68 (2), 614–625.
- Hollard, Guillaume, Sébastien Massoni, and Jean Christophe Vergnaud**, “In search of good probability assessors: an experimental comparison of elicitation rules for confidence judgments,” *Theory and Decision*, 2016, 80 (3), 363–387.
- Houssami, Nehmat, Les Irwig, Judy M Simpson, Merran McKessar, Steven Blome, and Jennie Noakes**, “The influence of clinical information on the accuracy of diagnostic mammography,” *Breast cancer research and treatment*, 2004, 85 (3), 223–228.
- Islam, Marco**, *Motivated Risk Assessments*, Department of Economics, School of Economics and Management, Lund University, 2021.

- Jevons, William S.**, *The Theory of Political Economy*, Macmillan and Company, 1879.
- Kappes, Andreas and Tali Sharot**, “The automatic nature of motivated belief updating,” *Behavioural Public Policy*, 2019, 3 (1), 87–103.
- Kőszegi, Botond**, “Emotional Agency,” *The Quarterly Journal of Economics*, 2006, 121 (1), 121–155.
- , “Utility from anticipation and personal equilibrium,” *Economic Theory*, 2010, pp. 415–444.
- Kőszegi, Botond and Matthew Rabin**, “Reference-dependent consumption plans,” *American Economic Review*, 2009, 99 (3), 909–36.
- Krizan, Zlatan and Paul D Windschitl**, “The influence of outcome desirability on optimism.,” *Psychological bulletin*, 2007, 133 (1), 95.
- Le Yaouanq, Yves**, “A model of ideological thinking,” *Mimeo, Ludwig Maximilians Universität München*, 2016, pp. 1–49.
- Lench, Heather C and Peter H Ditto**, “Automatic optimism: Biased use of base rate information for positive and negative events,” *Journal of Experimental Social Psychology*, 2008, 44 (3), 631–639.
- Leong, Yuan Chang, Brent L Hughes, Yiyu Wang, and Jamil Zaki**, “Neurocomputational mechanisms underlying motivated seeing,” *Nature Human Behaviour*, 2019.
- Lerman, Caryn, Chanita Hughes, Stephen J. Lemon, David Main, Carrie Snyder, Carolyn Durham, Steven Narod, and Henry T. Lynch**, “What you don’t know can hurt you: Adverse psychologic effects in members of BRCA1-linked and BRCA2-linked families who decline genetic testing,” *Journal of Clinical Oncology*, 1998, 16 (5), 1650–1654.
- Levy, Raphaël**, “Soothing politics,” *Journal of Public Economics*, 2014, 120, 126–133.
- Lim, Lena**, “A two-factor model of defensive pessimism and its relations with achievement motives,” *The Journal of Psychology*, 2009, 143 (3), 318–336.
- Loewenstein, George**, “Anticipation and the Valuation of Delayed Consumption,” *The Economic Journal*, 1987, 97 (387), 666–684.
- , “The Pleasures and Pains of Information,” *Science*, 2006, 312 (May), 704–706.

- **and Andras Molnar**, “The renaissance of belief-based utility in economics,” *Nature Human Behaviour*, 2018, 2 (3), 166–167.
- Macera, Rosario**, “Dynamic beliefs,” *Games and Economic Behavior*, 2014, 87, 1–18.
- Mayraz, Guy**, “Wishful Thinking,” *CEP Discussion Paper*, 2011, (1092).
- Mijović-Prelec, Danica and Drazen Prelec**, “Self-deception as self-signalling : a model and experimental evidence,” *Phil. Trans. of the Royal Society B*, jan 2010, 365 (1538), 227–240.
- Möbius, Markus M., Muriel Niederle, Paul Niehaus, and Tanya S. Rosenblat**, “Managing Self-Confidence,” *Mimeo, Stanford University*, 2014.
- Molnar, Andras and George Loewenstein**, “Thoughts and players: An Introduction to old and new economic perspectives on beliefs,” *The Science of Beliefs: A multidisciplinary Approach (provisional title, to be published in October 2021)*. Cambridge University Press. Edited by Julien Musolino, Joseph Sommer, and Pernille Hemmer, 2021.
- Mughan, A., C. Bean, and I. McAllister**, “Economic globalization, job insecurity and the populist reaction,” *Electoral Studies*, 2003, 22 (4), 617–633.
- Norem, Julie**, *The positive power of negative thinking*, Basic Books, 2008.
- Norem, Julie K and Nancy Cantor**, “Defensive pessimism: harnessing anxiety as motivation.,” *Journal of personality and social psychology*, 1986, 51 (6), 1208.
- Obschonka, Martin, Michael Stuetzer, Peter J. Rentfrow, Neil Lee, Jeff Potter, and Samuel D. Gosling**, “Fear, Populism, and the Geopolitical Landscape: The “Sleeper Effect” of Neurotic Personality Traits on Regional Voting Behavior in the 2016 Brexit and Trump Elections,” *Social Psychological and Personality Science*, 2018, 9 (3), 285–298.
- Orhun, A Yesim, Alain Cohn, and Collin Raymond**, “Motivated Optimism and Workplace Risk,” *Available at SSRN*, 2021.
- Oster, Emily, Ira Shoulson, and E. Ray Dorsey**, “Optimal expectations and limited medical testing: Evidence from huntington disease,” *American Economic Review*, 2013, 103 (2), 804–830.
- Pronk, Thomas, Dylan Molenaar, Reinout W Wiers, and Jaap Murre**, “Methods to split cognitive task data for estimating split-half reliability: A comprehensive review and systematic assessment,” *Psychonomic Bulletin & Review*, 2021, pp. 1–11.

- Salvador, Alexandre, Luc H Arnal, Fabien Vinckier, Philippe Domenech, Raphaël Gaillard, and Valentin Wyart**, “Premature commitment to uncertain decisions during human NMDA receptor hypofunction,” *Nature Communications*, 2022, *13* (1), 1–15.
- Santos-Pinto, Luís, Michele Dell, and Luca David Opromolla**, “A General Equilibrium Theory of Occupational Choice under Optimistic Expectations,” *Mimeo, University of Lausanne*, 2018.
- Schlag, Karl H., James Tremewan, and Joël J. Van der Weele**, “A Penny for Your Thoughts: A Survey of Methods for Eliciting Beliefs,” *Experimental Economics*, 2015, *18* (3), 457–490.
- Schmitz, Anja and Christian Grillon**, “Assessing fear and anxiety in humans using the threat of predictable and unpredictable aversive events (the NPU-threat test),” *Nature Protocols*, 2012, *7* (3), 527–532.
- Schwardmann, Peter**, “Motivated health risk denial and the market for preventative health care,” *Journal of Health Economics*, 2019, *In press*.
- , **Egon Tripodi, and Joël J van der Weele**, “Self-Persuasion: Evidence from Field Experiments at International Debating Competitions,” *American Economic Review*, 2022, *112* (4), 1118–46.
- Shah, Punit, Adam J. L. Harris, Geoffrey Bird, Caroline Catmur, and Ulrike Hahn**, “A pessimistic view of optimistic belief updating,” *Cognitive Psychology*, 2016, *90*, 71–127.
- Sharot, Tali, Alison M Riccardi, Candace M Raio, and Elizabeth A Phelps**, “Neural mechanisms mediating optimism bias,” *Nature*, 2007, *450* (7166), 102–105.
- , **Christoph W. Korn, and Raymond J. Dolan**, “How unrealistic optimism is maintained in the face of reality,” *Nature Neuroscience*, 2011, *14* (11), 1475–1479.
- , **Marc Guitart-Masip, Christoph W Korn, Rumana Chowdhury, and Raymond J Dolan**, “How dopamine enhances an optimism bias in humans,” *Current Biology*, 2012, *22* (16), 1477–1481.
- Simmons, Joseph P and Cade Massey**, “Is optimism real?,” *Journal of Experimental Psychology: General*, 2012, *141* (4), 630.
- Sinding Bentzen, Jeanet**, “Acts of God? Religiosity and Natural Disasters Across Subnational World Districts,” *Economic Journal*, 2019, *In press*.

- Sloman, Steven A., Philip M. Fernbach, and York Hagmayer**, “Self-deception requires vagueness,” *Cognition*, 2010, *115* (2), 268–281.
- Trautmann, Stefan T. and Gijs van de Kuilen**, “Belief Elicitation: A Horse Race among Truth Serums,” *Economic Journal*, 2014, *125* (589), 2116–2135.
- Weinstein, Neil D**, “Unrealistic optimism about susceptibility to health problems,” *Journal of behavioral medicine*, 1982, *5* (4), 441–460.
- Windschitl, Paul D, Andrew R Smith, Jason P Rose, and Zlatan Krizan**, “The desirability bias in predictions: Going optimistic without leaving realism,” *Organizational behavior and human decision processes*, 2010, *111* (1), 33–47.
- Wyart, Valentin and Etienne Koechlin**, “Choice variability and suboptimality in uncertain environments,” *Current Opinion in Behavioral Sciences*, 2016, *11*, 109–115.
- Yariv, Leeat**, “I’ll See It When I Believe It - A Simple Model of Cognitive Consistency,” *Available at SSRN 300696*, 2002.
- Zimmermann, Florian**, “The dynamics of motivated beliefs,” *American Economic Review*, 2020, *110* (2), 337–61.

# Appendices

## A Additional Tables and Figures

Table A.1: Average fraction of correct answers by experiment and treatment.

	Experiment 1 $N = 60$	Experiment 2 $N = 221$	Experiment 3 $N = 426$	Experiment 4 $N = 409$
Aggregate	70.5 (4.55)	68.5 (13.0)	75.6 (11.6)	75.3 (10.1)
No Shock/loss pattern	72.3 (6.45)	77.1 (16.9)	77.8 (12.9)	79.7 (11.9)
Shock/loss pattern	68.6 (6.57)	60.3 (17.5)	73.4 (15.3)	71.1 (16.4)
Difficulty Level 1 (easiest)	79.1 (6.17)	76.1 (17.6)	continuous	85.1 (12.8)
Difficulty Level 2	70.5 (6.08)	60.9 (11.6)	continuous	79.7 (12.9)
Difficulty Level 3	62.9 (7.95)	- -	continuous	72.6 (14.6)
Difficulty Level 4 (hardest)	- -	- -	continuous	64.1 (15.0)
Accuracy bonus Low	70.1 (5.68)	68.7 (14.2)	75.2 (12.6)	74.7 (11.4)
Accuracy bonus High	70.9 (5.59)	68.2 (14.7)	75.9 (75.6)	76.3 (12.6)
Low Stake	- -	68.5 (14.0)	- -	- -
High Stake	- -	68.7 (15.0)	- -	- -
No Stake	- -	71.2 (15.1)	- -	75.7 (10.4)

An observation is one individual's average accuracy in the specified condition. Standard deviations in brackets. Averages for Experiment 2 and 4 exclude the Neutral condition, which does not constitute a test of wishful thinking, and is reported separately, at the bottom.



## A.1 Experiment 1

Table A.2: Accuracy levels in Experiment 1

	(1)	(2)	(3)	(4)	(5)	(6)
	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy
Shock pattern	-3.698*** (1.202)	-3.073 (1.891)	-3.698*** (1.202)	-3.073 (1.891)	-4.111*** (1.264)	-2.833 (1.948)
High accuracy bonus (HAB)	0.885 (0.857)	0.174 (1.270)	0.885 (0.857)	0.174 (1.270)	0.785 (0.878)	0.313 (1.389)
Medium difficulty (MD)	-8.606*** (0.813)	-8.551*** (1.138)	-8.611*** (0.814)	-8.555*** (1.138)	-8.639*** (0.822)	-8.583*** (1.144)
High Difficulty (HD)	-17.15*** (1.270)	-14.60*** (1.584)	-17.14*** (1.272)	-14.59*** (1.585)	-17.20*** (1.269)	-14.64*** (1.593)
Shock pattern x HAB		1.424 (1.695)		1.424 (1.695)		0.944 (1.789)
Shock pattern x MD		-0.111 (1.889)		-0.111 (1.889)		-0.111 (1.895)
Shock pattern x HD		-5.097** (2.186)		-5.097** (2.186)		-5.139** (2.206)
Constant	80.50*** (1.020)	80.19*** (1.250)	80.50*** (0.871)	80.19*** (1.085)	80.76*** (1.056)	80.12*** (1.314)
Observations	11520	11520	11520	11520	720	720
$R^2$	0.019	0.020	0.020	0.020	0.261	0.268
ID clustering	✓	✓	✓	✓	✓	✓
Individual fixed effects			✓	✓		
Averages by ID/treatment					✓	✓

Linear regressions of accuracy on treatment dummies. Columns 1-4 are panel data regressions, columns 3-4 include individual fixed effects. Columns 5-6 are OLS regressions where each observation is average accuracy per treatment and individual. Standard errors in parentheses clustered by individual. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

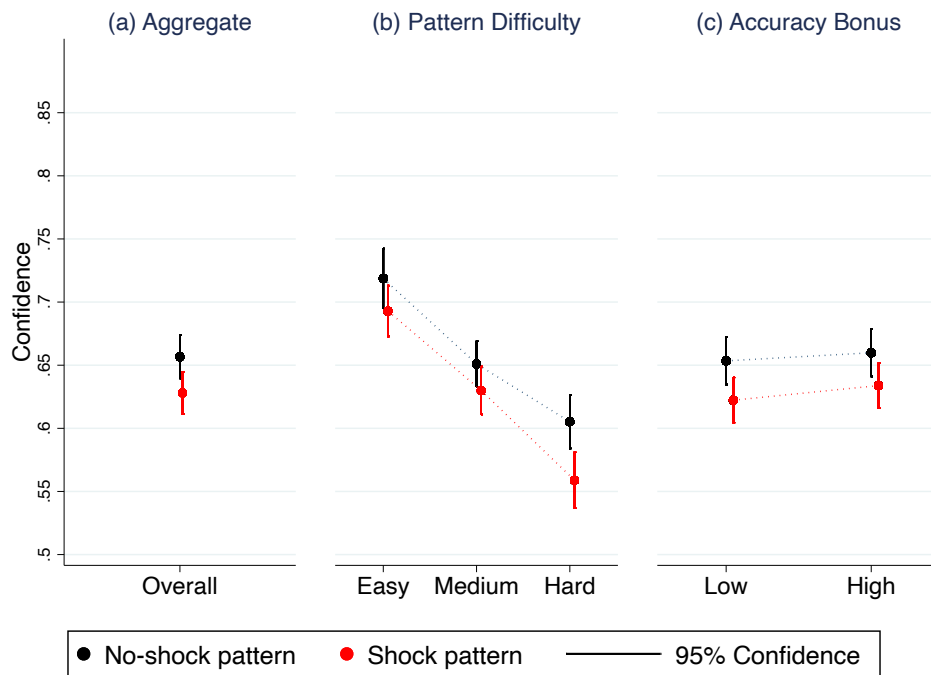


Figure A.1: **Electric shocks and confidence in the correct answer.** Average confidence levels in the correct answer, split by shock and no-shock pattern. Bars indicate 95% confidence intervals. One observation is the average over an individual's trials in a given category, so  $n = 60$  in each category. Panel a) shows aggregate results. Panel b) disaggregates the results by difficulty (tilt) of the pattern. Panel c) disaggregates by incentives for accuracy.

Table A.3: Belief levels in Experiment 1

	(1)	(2)	(3)	(4)	(5)	(6)
	Belief	Belief	Belief	Belief	Belief	Belief
Shock pattern	-2.858*** (0.938)	-2.835** (1.231)	-2.858*** (0.938)	-2.835** (1.231)	-3.108*** (0.977)	-2.943** (1.236)
High accuracy bonus (HAB)	0.898* (0.533)	0.642 (0.788)	0.898* (0.533)	0.642 (0.788)	0.841 (0.562)	0.478 (0.842)
Medium difficulty (MD)	-6.503*** (0.568)	-6.739*** (0.797)	-6.516*** (0.568)	-6.752*** (0.798)	-6.532*** (0.570)	-6.767*** (0.801)
High Difficulty (HD)	-12.34*** (0.959)	-11.32*** (1.224)	-12.32*** (0.965)	-11.30*** (1.228)	-12.37*** (0.956)	-11.34*** (1.228)
Shock pattern x HAB		0.510 (0.930)		0.510 (0.930)		0.726 (0.939)
Shock pattern x MD		0.472 (1.109)		0.472 (1.109)		0.470 (1.112)
Shock pattern x HD		-2.040 (1.348)		-2.040 (1.348)		-2.056 (1.360)
Constant	71.56*** (1.072)	71.55*** (1.231)	71.56*** (0.688)	71.55*** (0.827)	71.72*** (1.101)	71.63*** (1.251)
ID clustering	✓	✓	✓	✓	✓	✓
Individual fixed effects			✓	✓		
Averages by ID/treatment					✓	✓

Linear regressions of beliefs on treatment dummies. Beliefs are constructed from confidence judgments on a scale from 0 to 100, where the latter is perfect confidence in the correct answer. Columns 1-4 are panel data regressions, columns 3-4 include individual fixed effects. Columns 5-6 are OLS regressions where each observation is average accuracy per treatment and individual and correspond to those in Table 2. Standard errors in parentheses clustered by individual. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

## A.2 Experiment 2

Table A.4: Accuracy levels in Experiment 2

	(1)	(2)	(3)	(4)	(5)	(6)
	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy
Loss pattern	-16.59*** (1.531)	-9.155*** (3.284)	-16.59*** (1.531)	-9.143*** (3.287)	-16.54*** (1.605)	-8.248** (3.489)
High accuracy bonus (HAB)	-0.0284 (0.806)	-0.119 (1.004)	-0.0574 (0.812)	-0.157 (1.012)	-0.588 (0.851)	-1.081 (1.089)
Difficult pattern (DP)	-15.29*** (0.968)	-10.82*** (1.019)	-15.29*** (0.970)	-10.82*** (1.020)	-15.68*** (1.019)	-11.04*** (1.114)
Loss size (LS)	0.0380 (0.893)	1.100 (1.212)	0.0929 (0.892)	1.163 (1.209)	-0.617 (0.906)	0.776 (1.245)
Loss pattern x HAB		0.221 (1.618)		0.242 (1.620)		0.994 (1.771)
Loss pattern x DP		-8.806*** (1.586)		-8.795*** (1.586)		-9.200*** (1.701)
Loss pattern x LS		-2.079 (1.813)		-2.096 (1.815)		-2.784 (1.869)
Constant	84.57*** (1.944)	80.77*** (2.210)	84.54*** (1.719)	80.73*** (2.068)	85.82*** (1.964)	81.65*** (2.310)
Observations	11396	11396	11396	11396	3415	3415
$R^2$			0.064	0.066	0.134	0.140
ID clustering	✓	✓	✓	✓	✓	✓
Individual fixed effects			✓	✓		
Averages by ID/treatment					✓	✓

Linear regressions of accuracy on treatment dummies. Columns 1-4 are panel data regressions, columns 5-6 include individual fixed effects. Columns 5-6 are OLS regressions where each observation is average accuracy per treatment and individual and correspond to those in Table 2. Standard errors in parentheses clustered by individual. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

### A.3 Experiment 3

Table A.5: Accuracy levels in Experiment 3

	(1)	(2)	(3)	(4)	(5)	(6)
	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy
Loss pattern	-4.415*** (0.765)	-3.592*** (0.862)	-4.407*** (0.765)	-3.584*** (0.862)	-4.266*** (0.766)	-3.052*** (0.865)
High accuracy bonus (HAB)	0.728 (0.471)	0.484 (0.589)	0.723 (0.473)	0.479 (0.590)	0.630 (0.474)	0.685 (0.601)
Difficult pattern (DP)	-20.76*** (0.668)	-19.68*** (0.795)	-20.80*** (0.668)	-19.72*** (0.797)	-20.55*** (0.668)	-19.39*** (0.794)
Loss pattern x HAB		0.478 (0.838)		0.479 (0.838)		-0.110 (0.881)
Loss pattern x DP		-2.137* (0.849)		-2.138* (0.849)		-2.317** (0.892)
Constant	87.79*** (0.773)	87.37*** (0.812)	87.98*** (0.546)	87.56*** (0.574)	87.66*** (0.791)	87.06*** (0.829)
Observations	33507	33507	33507	33507	3408	3408
$R^2$			0.064	0.064	0.236	0.236
ID clustering	✓	✓	✓	✓	✓	✓
Individual fixed effects			✓	✓		
Averages by ID/treatment					✓	✓

Linear regressions of accuracy on treatment dummies. Columns 1-4 are panel data regressions, columns 3-4 include individual fixed effects. Columns 5-6 are OLS regressions where each observation is average accuracy per treatment and individual and correspond to those in Table 2. Standard errors in parentheses clustered by individual. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

## A.4 Experiment 4

Table A.6: Accuracy levels in Experiment 4

	(1)	(2)	(3)	(4)	(5)	(6)
	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy
Loss pattern	-8.445*** (1.030)	-6.560*** (1.285)	-8.463*** (1.031)	-6.567*** (1.285)	-8.453*** (1.040)	-7.393*** (1.308)
High accuracy bonus (HAB)	1.676*** (0.566)	1.175 (0.797)	1.738*** (0.571)	1.249 (0.803)	1.670*** (0.627)	1.004 (0.853)
Difficult pattern (DP)	-7.028*** (0.255)	-6.227*** (0.341)	-7.030*** (0.254)	-6.230*** (0.341)	-7.073*** (0.269)	-6.497*** (0.361)
Loss pattern x HAB		0.988 (1.200)		0.965 (1.201)		1.330 (1.322)
Loss pattern x DP		-1.578*** (0.468)		-1.577*** (0.468)		-1.150** (0.503)
Constant	89.37*** (0.704)	88.41*** (0.789)	89.31*** (0.671)	88.35*** (0.754)	89.53*** (0.732)	89.00*** (0.796)
Observations	21220	21220	21220	21220	6534	6534
$R^2$			0.046	0.046	0.110	0.110
ID clustering	✓	✓	✓	✓	✓	✓
Individual fixed effects			✓	✓		
Averages by ID/treatment					✓	✓

Linear regressions of accuracy on treatment dummies. Columns 1-4 are panel data regressions, columns 3-4 include individual fixed effects. Columns 5-6 are OLS regressions where each observation is average accuracy per treatment and individual and correspond to those in Table 2. Standard errors in parentheses clustered by individual. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

Table A.7: Accuracy levels in Experiment 4 split by dot counting behavior

	(1)	(2)	(3)	(4)
	Accuracy	Accuracy	Accuracy	Accuracy
Loss pattern	-8.034*** (1.792)	-11.80*** (2.075)	-2.619 (2.981)	-10.29*** (1.790)
High accuracy bonus (HAB)	0.732 (1.222)	0.534 (1.378)	5.238* (2.580)	1.317 (1.188)
Loss pattern x HAB	-0.192 (2.007)	4.038** (1.991)	-0.595 (3.090)	2.977* (1.689)
Constant	76.93*** (1.045)	80.71*** (1.151)	87.56*** (2.441)	81.77*** (1.033)
Counting dots:	Never	Sometimes	Always	Sometimes/Always
Observations	3400	2462	560	3134
$R^2$	0.022	0.035	0.019	0.030

OLS regressions of accuracy on treatment dummies, split by dot-counting behavior. Each observation is the average accuracy per treatment and individual. Standard errors in parenthesis are clustered at the individual level. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

Table A.8: Response times in Experiment 4 split by dot counting behavior

	(1)	(2)	(3)	(4)
	Response time	Response time	Response time	Response time
High accuracy bonus (HAB)	71.90 (47.68)	1001.0*** (242.8)	2285.6* (1219.3)	1384.7*** (332.8)
Constant	2470.5*** (98.01)	3358.3*** (212.2)	14189.1*** (1554.6)	4048.6*** (328.8)
Counting dots:	Never	Sometimes	Always	Sometimes&Always
Observations	3400	2462	560	3134
$R^2$	0.002	0.024	0.021	0.016

Quantile regression of median accuracy on treatment dummies, split by dot-counting behavior in different columns. Standard errors in parenthesis are clustered at the individual level. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

## A.5 Robustness



Table A.9: Accuracy and treatment effect in selected samples

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy
Shock pattern	-7.234*** (0.737)	-5.949*** (0.852)	-7.470*** (0.863)	-6.880*** (0.825)	-6.123*** (0.698)	-6.376*** (0.821)	-6.961*** (0.820)
High accuracy bonus (HAB)	0.672 (0.507)	0.994 (0.633)	0.526 (0.601)	0.946 (0.594)	0.713 (0.535)	0.627 (0.576)	1.308** (0.594)
Difficult pattern (DP)	-8.011*** (0.346)	-8.305*** (0.416)	-8.321*** (0.416)	-7.922*** (0.382)	-8.326*** (0.345)	-7.937*** (0.393)	-8.022*** (0.384)
Shock pattern x HAB	0.730 (0.782)	0.682 (0.969)	1.317 (0.944)	0.320 (0.933)	1.157 (0.821)	1.004 (0.873)	0.00287 (0.907)
Shock pattern x DP	-1.299*** (0.480)	-1.505*** (0.572)	-1.139* (0.588)	-1.618*** (0.532)	-1.453*** (0.466)	-1.294** (0.523)	-1.167** (0.534)
Constant	82.24*** (0.740)	77.25*** (1.228)	77.15*** (1.185)	76.35*** (1.109)	82.44*** (0.666)	81.56*** (0.845)	75.81*** (1.100)
Restrictions	None	Difficult Instructions < 4 out of 7	Hard to recall conditions < 5 out of 7	Total no. of mistakes on control q. < 3	Average accuracy > 60 percent	Trust in experimenters > 2 out of 5	Answer does not cause shock
ID clustered S.E.	✓	✓	✓	✓	✓	✓	✓
Experiment fixed effects	✓	✓	✓	✓	✓	✓	✓
Observations	12429	8372	8876	9742	11069	9647	9701
$R^2$	0.134	0.138	0.132	0.135	0.154	0.131	0.132

OLS regressions of accuracy on treatments across experiments. An observation is the average accuracy per treatment and individual. All regressions include experiment fixed effects and standard errors clustered at the individual level. “Shock pattern” is a dummy representing the pattern is associated with a shock (Experiment 1) or loss (Experiments 2-4). “High accuracy bonus” is a dummy that represents a high accuracy bonus, while “Difficult pattern” is a categorical variable that measures the difficulty of the perceptual task, (see 1 for exact specification by experiment). Column 2 excludes participants with one of the three highest scores on the question “How difficult did you find it to follow the instructions of this experiment?” measured on a 7-point Likert scale. Column 3 excludes participants with one of the three highest scores on the question “How difficult did you find it to keep in mind information about the potential losses and bonuses associated with trials” measured on a 7-point Likert scale. Column 4 excludes participants who wrongly answered more than two control questions at the beginning of the experiment. Column 5 excludes participants whose average accuracy in the experiment is below 60 percent. Column 6 excludes participants with the three lowest agreement scores with the statement “During the experiment, I never thought that I was deceived by the experimenters about my possible gains or losses” measured on a 5 point Likert scale. Column 7 excludes participants who wrongly answered a multiple choice question about the determinants of the shock. Data in column 2, 3, 4 and 7 exclude Experiment 1, where the relevant measure was not collected. \*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

Table A.10: Accuracy and treatment effect in selected samples, panel regressions

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy
Shock pattern	-7.234*** (0.737)	-5.949*** (0.852)	-6.987*** (0.938)	-6.880*** (0.825)	-6.123*** (0.698)	-6.376*** (0.821)	-6.961*** (0.820)
High accuracy bonus (HAB)	0.672 (0.507)	0.994 (0.633)	0.840 (0.687)	0.946 (0.594)	0.713 (0.535)	0.627 (0.576)	1.308** (0.594)
Difficult pattern (DP)	-8.011*** (0.346)	-8.305*** (0.416)	-8.236*** (0.453)	-7.922*** (0.382)	-8.326*** (0.345)	-7.937*** (0.393)	-8.022*** (0.384)
Shock pattern x HAB	0.730 (0.782)	0.682 (0.969)	1.381 (1.069)	0.320 (0.933)	1.157 (0.821)	1.004 (0.873)	0.00287 (0.907)
Shock pattern x DP	-1.299*** (0.480)	-1.505*** (0.572)	-1.429** (0.635)	-1.618*** (0.532)	-1.453*** (0.466)	-1.294** (0.523)	-1.167** (0.534)
Constant	82.24*** (0.740)	77.25*** (1.228)	77.13*** (1.329)	76.35*** (1.109)	82.44*** (0.666)	81.56*** (0.845)	75.81*** (1.100)
Restrictions	None	Difficult Instructions < 4 out of 7	Hard to recall conditions < 4 out of 7	Total no. of mistakes on control q. < 3	Average accuracy > 60 percent	Trust in experimenters > 2 out of 5	Answer does not cause shock
ID clustered S.E.	✓	✓	✓	✓	✓	✓	✓
Individual fixed effects	✓	✓	✓	✓	✓	✓	✓
Observations	12429	8372	7372	9742	11069	9647	9701
$R^2$	0.134	0.138	0.136	0.135	0.154	0.131	0.132

Linear regressions of accuracy on treatments across experiments. We use a panel data structure where each observation is a single trial, and regressions include individual fixed effects and standard errors clustered at the individual level. “Shock pattern” is a dummy if the pattern is associated with a shock (Experiment 1) or loss (Experiments 2-4). “High accuracy bonus” is a dummy that represents a high accuracy bonus, while “Difficult pattern” is a categorical variable that measures the difficulty of the perceptual task, (see 1 for exact specification by experiment). Column 2 excludes participants with one of the three highest scores on the question “How difficult did you find it to follow the instructions of this experiment?” measured on a 7-point Likert scale. Column 3 excludes participants with one of the three highest scores on the question “How difficult did you find it to keep in mind information about the potential losses and bonuses associated with trials” measured on a 7-point Likert scale. Column 4 excludes participants who wrongly answered more than two control questions at the beginning of the experiment. Column 5 excludes participants whose average accuracy in the experiment is below 60 percent. Column 6 excludes participants with the three lowest agreement scores with the statement “During the experiment, I never thought that I was deceived by the experimenters about my possible gains or losses” measured on a 5 point Likert scale. Column 7 excludes participants who wrongly answered a multiple choice question about the determinants of the shock. Data in column 2, 3, 4 and 7 exclude Experiment 1, where the relevant measure was not collected. \*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

## A.6 Noise-based explanation and neutral patterns

Table A.11: Comparison of accuracy of neutral, loss and no-loss patterns

	(1)	(2)
	Accuracy	Accuracy
Loss pattern	-0.110*** (0.0124)	-0.0452*** (0.00823)
No-loss pattern	0.0586*** (0.00899)	0.0396*** (0.00626)
Constant	0.712*** (0.0102)	0.757*** (0.00517)
	Experiment 2	Experiment 4
Observations	1326	1227
$R^2$	0.123	0.065

OLS regression of accuracy on neutral, loss and no-loss patterns in Experiment 2 and 4. Baseline are neutral patterns where no shock was administered present. Each observation is average accuracy per treatment and individual. Standard errors clustered at the participant level in parentheses. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

## A.7 Dynamics

Table A.12: The dynamics of wishful thinking

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy
Loss pattern	-0.0652*** (0.0209)	-0.172*** (0.0207)	-0.0244** (0.0100)	-0.0747*** (0.0186)	-0.0559*** (0.0167)	-0.170*** (0.0197)	-0.0359*** (0.00784)	-0.0637*** (0.0166)
Trial	0.000411*** (0.000138)	0.000464 (0.000299)	0.000164 (0.000158)	0.000251 (0.000266)				
Loss pattern x Trial	0.000260 (0.000156)	0.000159 (0.000458)	-0.000477** (0.000233)	-0.000229 (0.000387)				
Second half					0.0519*** (0.0181)	0.0295* (0.0154)	0.000352 (0.00704)	0.0290** (0.0117)
Loss pattern x Second half					0.0342 (0.0221)	0.00731 (0.0242)	-0.0157* (0.00937)	-0.0279 (0.0180)
Experiment	1	2	3	4	1	2	3	4
Observations	11520	11396	33507	21220	11520	11396	33507	21220
$R^2$	0.008	0.035	0.003	0.010	0.008	0.035	0.003	0.011

OLS regression of accuracy on shock/loss and no-shock/no-loss patterns and temporal indices in various experiments. All regressions include participant fixed effects. Standard errors in parenthesis clustered at the participant level. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

Table A.13: Accuracy and response to previous losses

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy
Loss pattern	-0.0356*** (0.0124)	-0.168*** (0.0157)	-0.0420*** (0.00795)	-0.0827*** (0.0108)	-0.0638*** (0.0216)	-0.175*** (0.0209)	-0.0222** (0.0101)	-0.0729*** (0.0188)
Lagged loss	-0.00462 (0.0170)	-0.0165 (0.0141)	-0.00731 (0.00821)	0.0117 (0.0108)	-0.00409 (0.0170)	-0.0168 (0.0141)	-0.00699 (0.00821)	0.0115 (0.0108)
Loss pattern x Lagged loss	-0.00933 (0.0268)	0.0173 (0.0226)	-0.0149 (0.0127)	-0.0100 (0.0171)	-0.00950 (0.0270)	0.0166 (0.0226)	-0.0157 (0.0127)	-0.0100 (0.0171)
Trial					0.000411*** (0.000138)	0.000464 (0.000299)	0.000166 (0.000158)	0.000247 (0.000266)
Loss pattern x Trial					0.000260 (0.000156)	0.000161 (0.000458)	-0.000484** (0.000232)	-0.000229 (0.000387)
Experiment	1	2	3	4	1	2	3	4
Observations	11520	11396	33507	21220	11520	11396	33507	21220
$R^2$	0.002	0.035	0.003	0.010	0.008	0.035	0.003	0.010

OLS regression of accuracy on shock/loss and no-shock/no-loss patterns and lagged shocks or losses. All regressions include participant fixed effects. Standard errors in parenthesis clustered at the participant level. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

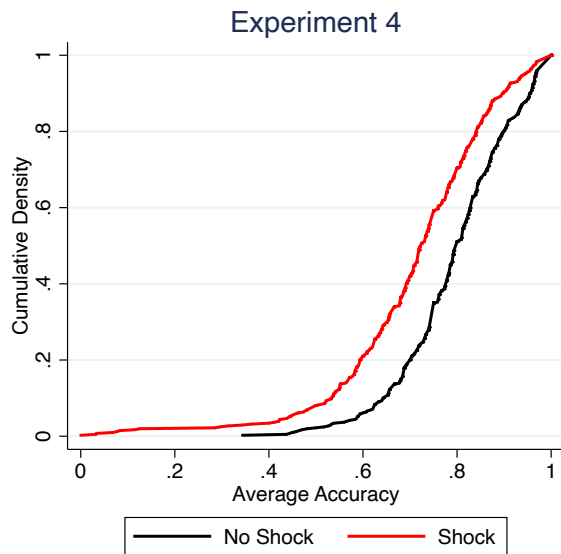
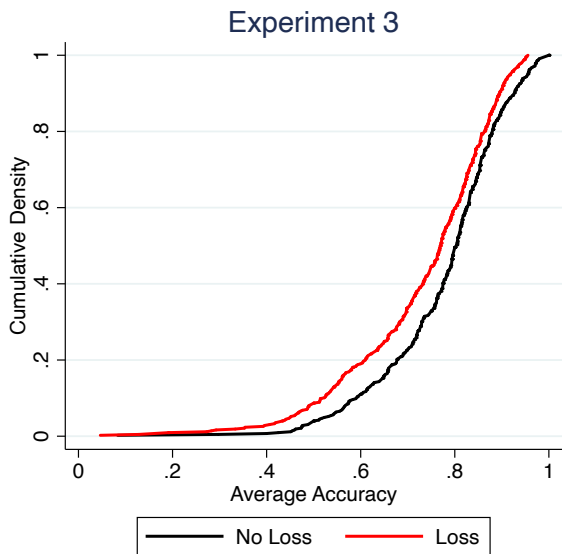
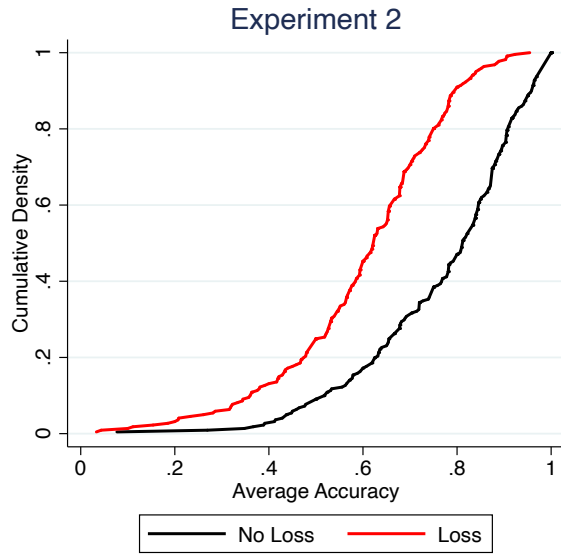
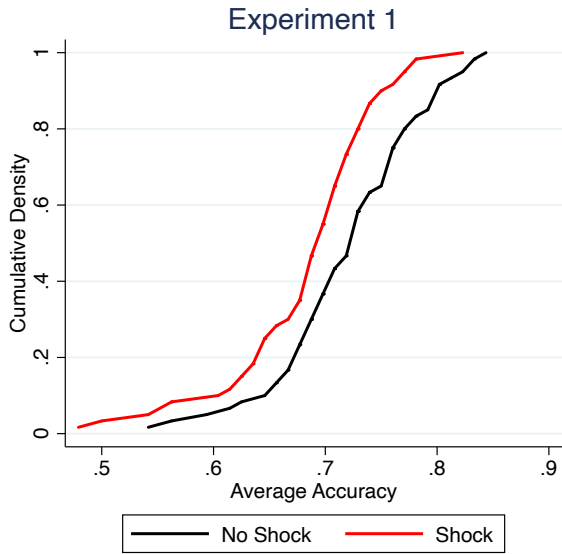


Figure A.2: CDFs of participants' average accuracy in each experiment, split by shock/loss and no-shock/no-loss patterns. Each observation is the average accuracy of all trials of an individual participant in that category.

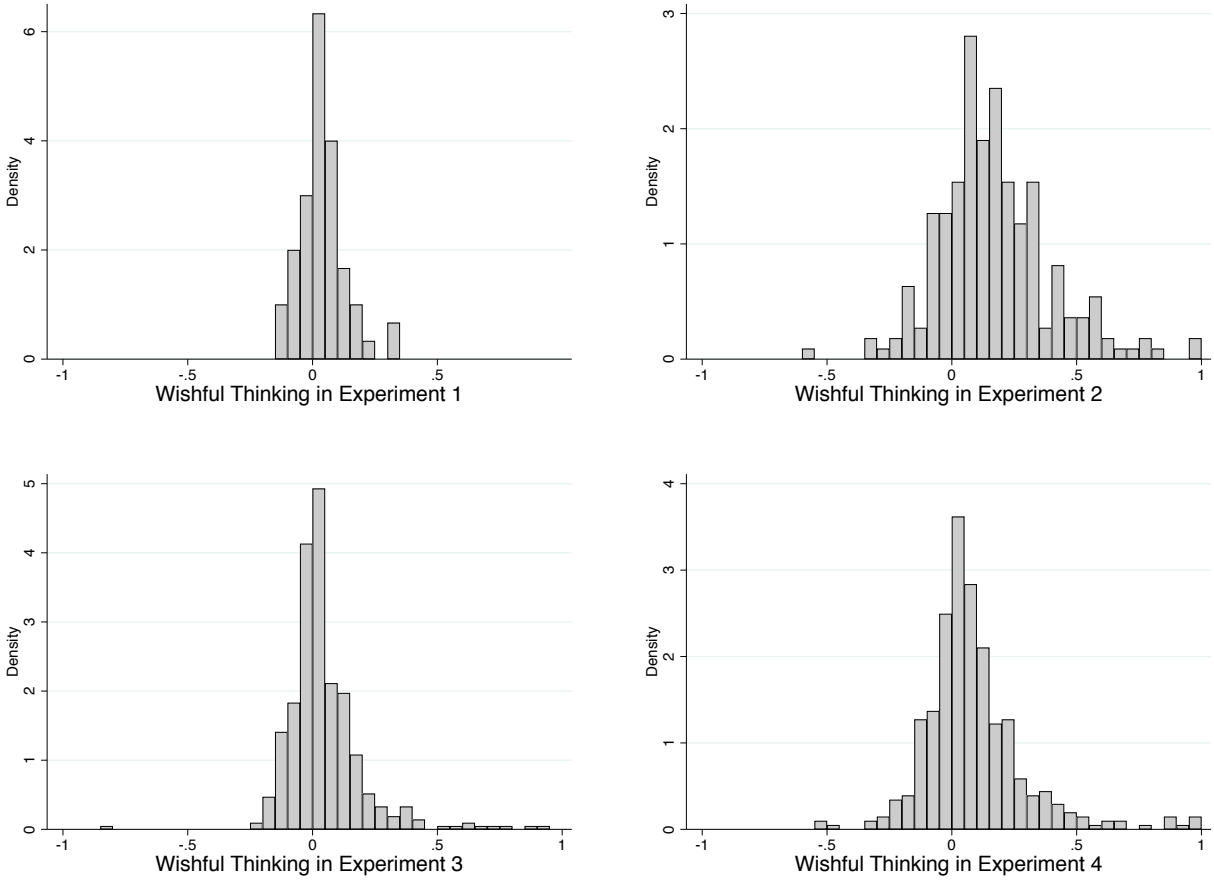
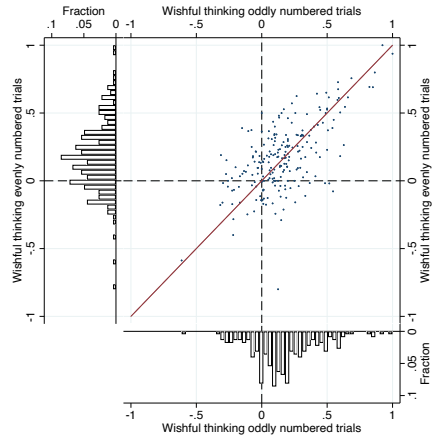
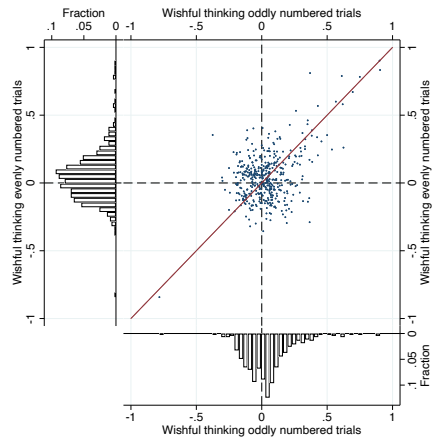


Figure A.3: Histograms of participants' wishful thinking in each experiment. Wishful thinking is defined as an individual's accuracy for shock patterns minus their accuracy for no shock patterns.

### A - Experiment 2



### B - Experiment 3



### C - Experiment 4

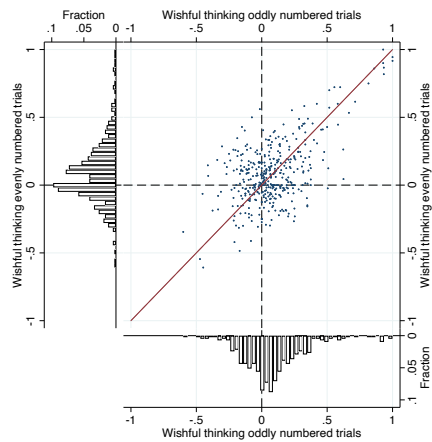
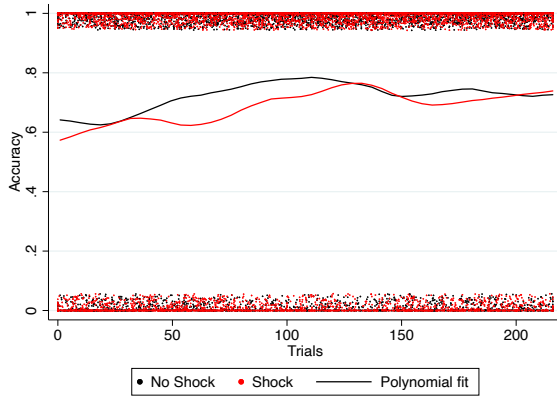


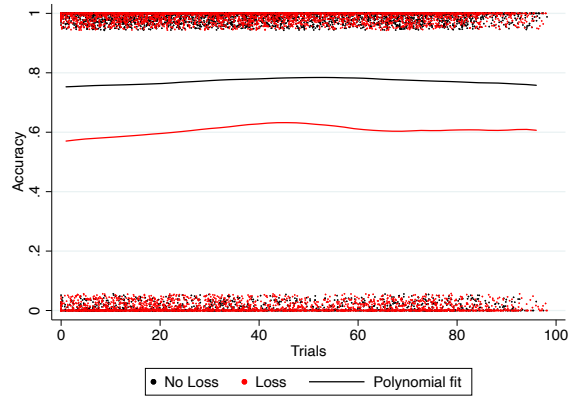
Figure A.4: Scatterplot of participants' wishful thinking in odd and even trials in each experiment. Wishful thinking is defined as an individual's accuracy for shock patterns minus their accuracy for no-shock patterns. Each plot includes the 45 degree line as well as projected histograms for odd and even trials.



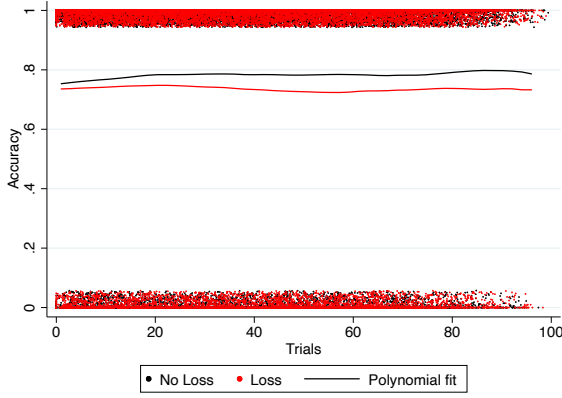
A - Experiment 1



B - Experiment 2



C - Experiment 3



D - Experiment 4

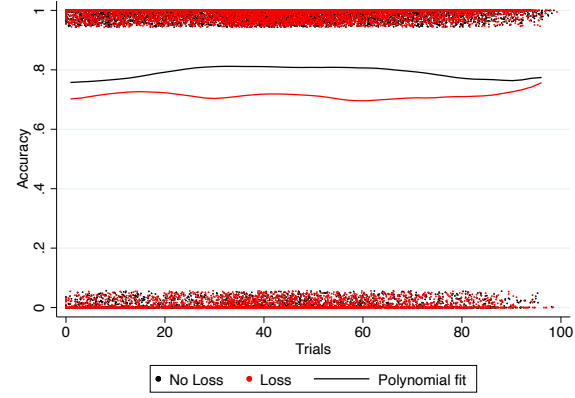


Figure A.5: Dynamics of wishful thinking over time by experiment. x-Axis shows trial numbers. y-Axis is accuracy in single trials (dots - jitter added) and a polynomial fit (line). Note that in Experiment 2, 3, and 4, not all participants completed the same number of trials, as trials stopped when participants reached 5 cumulative (and stochastic) losses from their endowment.

## B Replication Experiment 1

Before conducting Experiment 1, we ran an experiment with an almost identical design, which was also preregistered on [aspredicted.org](https://aspredicted.org) ( preregistration is here). There were some small differences. First, the experiment also featured some neutral trial blocks in which subjects did not face the threat of a shock. Second, while we used the same visual patterns for the task, subjects had to indicate whether they were vertically or horizontally oriented (rather than choosing the closest diagonal), and there were four difficulty levels. Third, incentives were constant across the experiment. Finally, the experimental code exhibited a small bug which meant that the ambiguity levels were not equally calibrated across the Shock and No-shock condition. While we are able to control for the ambiguity level (see below), the imperfect randomization ultimately caused us to rerun this experiment, resulting in Experiment 1. For our purposes, two aspects of the precursor experiment are of interest. First, we investigate whether our main treatment effect obtains also in this study. The experiment replicates our main results with very similar effect sizes ( $\mu$ ) between shock and no-shock patterns (Accuracy:  $\mu = 0.047$ ,  $s.d. = 0.103$ ,  $p = 0.0024$ ; Belief:  $\mu = 0.031$ ,  $s.d. = 0.069$ ,  $p = 0.0022$ ). Table B.1 shows the result of an OLS regression of Accuracy and Belief, averaged by subject and condition, on treatment dummies. Because pattern difficulty was not well balanced between the Shock and No-Shock conditions, we control for difficulty in Table B.1. This shows similar results, with a highly significant effect of the Shock condition which is 0.037 for Accuracy and 0.029 for Belief, closely mirroring the magnitudes in our main experiment.

Second, the presence of “Neutral” trials without the threat of shocks allows us to investigate the “noise explanation” elaborated in Section 5.5.3, where we inspect how the presence of a shock-threat affects Accuracy and Beliefs. We find that the Accuracy and Belief in the Neutral patterns are substantially worse than in the No-Shock patterns. This shows that the presence of a shock does not necessarily reduce performance, reinforcing our confidence that wishful thinking is driving our main result.

	(1)	(2)
	Accuracy	Belief
Shock Pattern	-0.0373** (0.0152)	-0.0285*** (0.0103)
Neutral Pattern	-0.0339*** (0.0114)	-0.0371*** (0.00736)
Pattern Difficulty	-0.0349*** (0.00497)	-0.0255*** (0.00411)
Constant	0.904*** (0.0132)	0.810*** (0.0127)
Observations	600	600
$R^2$	0.098	0.087

Table B.1: OLS regressions of Accuracy and Belief on the experimental conditions. Standard errors in parentheses clustered at subject level. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

## C Theory extensions

### C.1 Anticipatory utility from accuracy incentives

In the following, we suppose that the agent obtains utility from her anticipation of incentives for accuracy. Her utility then takes the following form.

$$\begin{aligned}
 U = & \frac{1}{2} (1 + 2p\hat{p} - \hat{p}^2) M + \sigma_m \frac{1}{2} (1 + \hat{p}^2) M \\
 & - (r_z p + (1 - r_z)(1 - p))qZ - \sigma_z (r_z \hat{p} + (1 - r_z)(1 - \hat{p}))qZ \\
 & - \lambda(s)(p - \hat{p})^2,
 \end{aligned}$$

where  $\sigma_m$  captures the agent's propensity to savor of future payoffs from the accuracy incentives and is distinct from the anxiety of being shocked or losing money, which is still parameterized by  $\sigma_z$ . Note that the agent's anticipatory utility from expecting future accuracy payoffs depends only on her chosen belief  $\hat{p}$  and not on the undistorted or true probability of the pattern being right-tilted  $p$ .

Maximizing the agent's utility then yields the following optimal beliefs

$$\hat{p}^* = \frac{M + 2\lambda(s)}{(1 - \sigma_m)M + 2\lambda(s)} p(s, r_t) - \frac{\sigma_z(2r_z - 1)qZ}{(1 - \sigma_m)M + 2\lambda(s)},$$

and wishful thinking is given by

$$W := \mathbb{E}_{r_t} [\hat{p}^*(r_z = 0) - \hat{p}^*(r_z = 1)] = \frac{2\sigma_z qZ}{(1 - \sigma_m)M + 2\lambda(s)}.$$

We see that an anticipatory utility motive stemming from the savoring of accuracy incentives does not change qualitative predictions regarding the drivers of anxiety-induced wishful thinking. In particular, the comparative statics of wishful thinking in  $M$ ,  $s$  and  $Z$  have the same sign they had in the main model. However, wishful thinking is now increasing in the savoring parameter  $\sigma_m$  because the higher  $\sigma_m$ , the more the agent cares about her perceived rather than her actual receipt of accuracy incentives, and the former is not decreasing in the amount of belief distortion.

## C.2 Defensive pessimism and bracing

To capture defensive pessimism or bracing we suppose that the agent maximizes the following utility function.<sup>26</sup>

$$\begin{aligned}
 U = & \frac{1}{2} (1 + 2p\hat{p} - \hat{p}^2) M \\
 & - (r_z p + (1 - r_z)(1 - p))q(Z - b\hat{p}Z) - \sigma_z(r_z\hat{p} + (1 - r_z)(1 - \hat{p}))q(Z - b\hat{p}Z) \\
 & - \lambda(s)(p - \hat{p})^2.
 \end{aligned}$$

Here, parameter  $b$  captures the benefit of bracing or the extent to which defensive pessimism can soften the blow of a shock or loss. If  $b = 0$ , pessimistic beliefs do not lessen the impact of a shock. In the other extreme,  $b = 1$ , the agent can fully negate the shock's impact by being maximally pessimistic. The agent's optimal beliefs are then given by

$$\begin{aligned}
 \hat{p}^*(r_z = 1) &= \frac{M + 2\lambda(s) + bZ}{M + 2\lambda(s) - 2\sigma_z qbZ} p - \frac{\sigma_z qZ}{M + 2\lambda(s) - 2\sigma_z qbZ} \\
 \hat{p}^*(r_z = 0) &= \frac{M + 2\lambda(s) - bZ}{M + 2\lambda(s) + 2\sigma_z qbZ} p + \frac{\sigma_z qZ}{M + 2\lambda(s) + 2\sigma_z qbZ}.
 \end{aligned}$$

We can check that  $\frac{\hat{p}^*(r_z=1)}{db} > 0$  and  $\frac{\hat{p}^*(r_z=0)}{db} < 0$ . So the bracing motive decreases apparent wishful thinking, which is defined as  $W := \hat{p}^*(r_z = 0) - \hat{p}^*(r_z = 1)$ . Furthermore, it is easy to verify that there is wishful thinking, i.e.  $\hat{p}^*(r_z = 1) < p$  and  $\hat{p}^*(r_z = 0) > p$ , only if the following inequality holds

$$b < \frac{\sigma_z q}{1 + 2\sigma_z q}.$$

This inequality is satisfied for small  $b$  and large  $\sigma_z$ . Therefore, if we detect wishful thinking in our average participant, then we provide evidence not just for wishful thinking but for the fact that the anticipatory anxiety motive dominates the defensive pessimism or bracing motive.

## C.3 The correlation between anxiety and wishful thinking

The agent's utility function contains a term that captures her experienced anxiety conditional on her (optimal) belief and the tilt of the shock pattern. Since the shock pattern is right-tilted 50

---

<sup>26</sup>We take "bracing" to mean any action or investment that can reduce the impact (physical or otherwise) of negative news or events. Defensive pessimism is a more specific concept and form of bracing. It is a cognitive strategy that allows one to deal with the psychological impact of negative events by holding negative expectations (Norem and Cantor, 1986).

percent of the time and left-tilted 50 percent of the time, average experienced anxiety is given by

$$A = \frac{1}{2}\sigma_z \hat{p}_{r_z=1}^* qZ + \frac{1}{2}\sigma_z (1 - \hat{p}_{r_z=0}^*) qZ$$

Substituting into this term her respective optimal beliefs  $\hat{p}_{r_z=1}^*$  and  $\hat{p}_{r_z=0}^*$  from the main model yields

$$A = \frac{1}{2}\sigma_z qZ - \frac{\sigma_z^2 q^2 Z^2}{M + 2\lambda}.$$

Next, we turn to the comparative statics of  $A$  in  $\lambda$  and  $\sigma_z$ . Comparing these with the comparative statics of  $W$  in  $\lambda$  and  $\sigma_z$ , will allow us to see how heterogeneities in  $\lambda$  and  $\sigma_z$  map into correlations between experienced anxiety and wishful thinking.

It is easy to see that  $\frac{dA}{d\lambda} > 0$ , so that experienced anxiety is increasing in the cognitive costs of self-deception. How  $A$  varies with  $\sigma_z$  is more subtle. We can show that  $\frac{dA}{d\sigma_z} > 0$  if and only if

$$\sigma_z < \frac{4q^2 Z^2}{M + 2\lambda}, \tag{C.1}$$

and  $\frac{dA}{d\sigma_z} \leq 0$  otherwise. So experienced anxiety  $A$  is increasing in innate anxiety  $\sigma_z$  for small levels of innate anxiety and decreasing for higher levels. To see why experienced anxiety must eventually decrease in innate anxiety note that for very high levels of innate anxiety the agent will engage in sufficient wishful thinking to put zero subjective probability on the negative outcome, leaving her with no experienced anxiety.

We can use our experimental results to get a sense of where in the parameter space participants are likely to be located. To this end, note that inequality C.1 can be rewritten to state that  $\frac{dA}{d\sigma_z} > 0$  if and only if

$$W < \frac{1}{2}$$

Clearly, average wishful thinking is well below 50 percentage points in all experiments. Going forward we therefore assume that inequality C.1 is satisfied.

Putting together the comparative statics on  $\sigma_z$  and  $\lambda$ , we now consider two scenarios.

**Secenario A (heterogeneity in  $\lambda$ ).** Suppose there are two groups of participants (group 1 and group 2) that differ only in  $\lambda$ , their cognitive costs of self-deception. In particular, suppose that  $\lambda_1 < \lambda_2$ , where the subscript is the group label. Based on how  $W$  and  $A$  vary with  $\lambda$ , it will then be the case that  $W_1 > W_2$  and  $A_1 < A_2$ . Therefore, in this scenario, wishful thinking and

experienced anxiety are negatively correlated, as those with low cost of self-deception have higher wishful thinking and hence lower experienced anxiety.

**Secenario B (heterogeneity in  $\sigma_z$ ).** Now suppose there are two groups of participants (group 1 and group 2) that differ only in  $\sigma_z$ , their innate anxiety. In particular, suppose that  $\sigma_{z1} < \sigma_{z2}$ , where the subscript is the group label. Based on how  $W$  and  $A$  vary with  $\sigma_z$ , it will then be the case that  $W_1 < W_2$  and  $A_1 < A_2$ . Therefore, in this scenario, wishful thinking and experienced anxiety are positively correlated, as high innate anxiety leads to the higher experienced anxiety, despite a partial offset by higher wishful thinking.

Considering these two scenarios allows us to conclude that, according to the model, whether measures of experienced anxiety and wishful thinking are positively, negatively, or not correlated is ambiguous. More specifically, it will depend on whether there is a dominant heterogeneity and whether this heterogeneity is in participants' ability to self-deceive  $\sigma_z$  (negative correlation) or in their innate anxiety  $\lambda$  (positive correlation).

## C.4 Optimal beliefs with a hard cognitive constraint

In this section we describe a model that better captures the statistical relationships we see in the data. We will add two elements to the model. First, self-deception is constrained to some maximum amount or hard cognitive constraint. One interpretation, closest to our original model and developed below in more detail, is that this is a binding constraint on an optimizing agent who chooses optimal beliefs. A second interpretation is that self-deception is an “automatic” or “system 1” process, where a certain amount of self-deception occurs automatically without an agent’s cognitive influence. The second new element is that the constraint on self-deception depends on the signal strength, which is determined by the agent’s investment in signal precision or information-gathering.

Thus, we suppose that the agent first invests in signal precision at time  $t = 0$  and then distorts her mental representation of the signal at time  $t = 1$ . To solve the agent’s problem, we first look at  $t = 1$ .

### C.4.1 Belief choice conditional on the signal at $t = 1$

Consider the following maximization problem.

$$\begin{aligned}
Max \quad U_1 &= \frac{1}{2} (1 + 2p\hat{p} - \hat{p}^2) M \\
&\quad - (r_z p + (1 - r_z)(1 - p))qZ - \sigma_Z(r_z\hat{p} + (1 - r_z)(1 - \hat{p}))qZ \\
\text{such that} \quad &|p - \hat{p}| \leq \lambda(s)
\end{aligned}$$

So self-deception is cognitively costless up to a point and then becomes impossible. In line with our results, we further assume that the maximum distance between true and distorted beliefs is decreasing in  $s$ , the signal precision, i.e.  $\lambda'(s) < 0$ .

Solving this maximization problem yields the following optimal beliefs.

$$\hat{p}^* = \begin{cases} \max(p - \frac{\sigma q Z}{M}, p - \lambda(s)) & \text{if } r_z = 1 \\ \min(p + \frac{\sigma q Z}{M}, p + \lambda(s)) & \text{if } r_z = 0 \end{cases}$$

Our results in Experiments 1 to 3 indicate that ex-post signal distortion responds to  $s$ , but not to  $M$  or  $Z$ . In other words, the cognitive constraint is binding for all values of  $M$  and  $Z$  and optimal beliefs are given by  $p - \lambda(s)$  for right-tilted patterns and by  $p + \lambda(s)$  for left-tilted patterns. In this case, wishful thinking is given by  $W = 2\lambda(s)$ , which is decreasing in signal precision  $s$ .

In a next step we look the an agent's investment in signal precision.

#### C.4.2 Investment in signal precision at $t = 0$

The agent decides on her investment in signal precision knowing whether shocks are associated with left or with right-tilted patterns and anticipating the effect of signal precision on her payments from the BDM mechanism and her ability to self-deceive. At the point of deciding on the cognitive effort she spends on identifying the pattern, she does not know the actual tilt of the pattern and merely has a 50:50 prior over whether the pattern is left- or right-tilted. Her choice of signal precision  $s$  maximizes the following function.

$$\begin{aligned}
U_0 &= \frac{1}{2} \left( \hat{p}_{r_t=1} M + (1 - \hat{p}_{r_t=1}) \frac{1 + \hat{p}_{r_t=1}}{2} M \right) + \frac{1}{2} \left( (1 - \hat{p}_{r_t=0}) \frac{1 + \hat{p}_{r_t=0}}{2} M \right) \\
&\quad - \frac{1}{2} \sigma_Z (r_z \hat{p}_{r_t=1} + (1 - r_z)(1 - \hat{p}_{r_t=1})) qZ - \frac{1}{2} \sigma_Z (r_z \hat{p}_{r_t=0} + (1 - r_z)(1 - \hat{p}_{r_t=0})) qZ \\
&\quad - \frac{1}{2} qZ - c(s),
\end{aligned} \tag{C.2}$$



where  $c(s)$  is the cognitive cost associated with generating more precise representations of the signal. Moreover,  $\hat{p}_{r_t=\{0,1\}}$  are the agent's optimal  $t = 1$  beliefs, conditional on the true pattern, that depend on  $p_{r_t=\{0,1\}}$ , her undistorted belief, which in turn depends on  $s \in [0.5, 1]$ .

For simplicity, we assume that undistorted beliefs  $p_{r_t=1} = s$  and  $p_{r_t=0} = 1 - s$ . So as  $s$  increases, the agent becomes more accurate in identifying both right- and left-tilted patterns. Furthermore, we assume that  $\lambda(s) = \epsilon(1 - s)$  and that  $c(s) = cs^2$ .

We consider the case of  $r_z = 1$  so that  $\hat{p}_{r_t} = p_{r_t} - \lambda(s)$ . Then, substituting the expressions for  $p_{r_t=\{0,1\}}$ ,  $\lambda(s)$  and  $c(s) = cs^2$  into C.2 and simplifying yields

$$U_0 = \frac{1}{2}M \left( \frac{1}{2} + 2s - s^2 - \frac{1}{2}\epsilon^2(1 - s) \right) - \frac{1}{2}\sigma_Z(1 - \epsilon + \epsilon s)qZ - \frac{1}{2}qZ - cs^2$$

The  $s$  that maximizes this ex ante utility is given by

$$s^* = \frac{\frac{1}{2}M(2 + \epsilon^2) - \frac{1}{2}\sigma_z\epsilon qZ}{\frac{1}{2}M(2 + \epsilon^2) + 2c} \quad (\text{C.3})$$

Taking the first derivative yields that  $\frac{ds^*}{dM} > 0$ . So an increase in accuracy incentives increases signal precision. Then, because  $\lambda'(s) < 0$  and wishful thinking is decreasing in  $\lambda$ , we have that  $\frac{dW}{dM} < 0$ .

Note that this chain of reasoning depends on the agent being *able* to increase her signal precision, i.e.  $\frac{ds^*}{dM} > 0$ . We can proxy signal precision by performance as exogenous changes in signal precision (or ambiguity) lead to clear variation in performance in all our experiment. By this is true for dot counters in experiment 4, but not for participants in experiments 1 to 3.

**Hypothesis C.1 (Incentives)** *If signal precision and hence, accuracy, is increasing in accuracy incentives  $\frac{ds^*}{dM} > 0$ , then wishful thinking acts decreasing in accuracy incentives, i.e.  $\frac{dW}{dM} < 0$ .*

**Naivite and sophistication.** Here we assumed that the agent is sophisticated about the effect of her investment in  $s$  on her ability to self-deceive. An agent who is naive about the link between signal precision and her subsequent ability to self-deceive expects that  $\lambda'(s) = 0$ . The naive agent's

optimal signal precision is then given by

$$s_n^* = \frac{M}{M + 2c} \tag{C.4}$$

For the naive it will therefore also be the case that  $\frac{ds_n^*}{dM} > 0$ . Because, in reality,  $\lambda(s) > 0$  and  $\lambda'(s) < 0$ , a naive's wishful thinking will then also be decreasing in  $M$ . This implies that the result that wishful thinking is decreasing in  $M$  does not help us distinguish between naives and sophisticates. However, note that  $s^*$  but not  $s_n^*$  are increasing in  $Z$  the size of the loss or shock. A positive effect of  $Z$  on wishful thinking is therefore suggestive of sophistication.

**Ethics Committee Economics and Business (EBEC)  
University of Amsterdam**

**Amsterdam Business School**  
Plantage Muidergracht 12  
1012 TV Amsterdam  
The Netherlands  
T +31 20 525 7384  
[www.abs.uva.nl](http://www.abs.uva.nl)

To: van der Weele

Date	Our reference	
June 12, 2017	EC 20170510120541	
Contact	Telephone	E-Mail
Sophia de Jong	(31)20-5255311	secbs-abs@uva.nl
Subject		
EBEC approval		

Dear Joel van der Weele,

The Economics & Business Ethics Committee (University of Amsterdam) received your request nr 20170510120541 to approve your project "Anticipatory utility and probabilistic confidence".

We evaluated your proposed research in terms of potential impact of the research on the participants, the level and types of information and explanation provided to the participants at various stages of the research process, the team's expertise in conducting the proposed analyses and particularly in terms of restricted access to the data to guarantee optimal levels of anonymity to the participants.

The Ethics Committee approves of your request.

Best regards,

On behalf of the Ethics Committee Economics and Business,

Prof. Dr. J.H. Sonnemans  
Chairman of the Committee

**Ethics Committee Economics and Business (EBEC)  
University of Amsterdam**

**Amsterdam Business School**  
Plantage Muidergracht 12  
1012 TV Amsterdam  
The Netherlands  
T +31 20 525 7384  
www.abs.uva.nl

To: van der Weele

Date	Our reference	
February 02, 2021	EC 20210202020244	
Contact	Telephone	E-Mail
Sophia de Jong	(31)20-5255311	secbs-abs@uva.nl
Subject		
EBEC approval		

Dear Joel van der Weele,

The Economics & Business Ethics Committee (University of Amsterdam) received your request nr 20210202020244 to approve your project "Anticipatory anxiety from monetary losses and wishful thinking".

We evaluated your proposed research in terms of potential impact of the research on the participants, the level and types of information and explanation provided to the participants at various stages of the research process, the team's expertise in conducting the proposed analyses and particularly in terms of restricted access to the data to guarantee optimal levels of anonymity to the participants.

The Ethics Committee approves of your request.

The information as filled in the form, can be found at  
<https://www.creedexperiment.nl/EBEC/showprojectAVG.php?nummer=20210202020244>

Best regards,

On behalf of the Ethics Committee Economics and Business,

Prof. Dr. J.H. Sonnemans  
Chairman of the Committee

## Wishful thinking and anxiety in the laboratory (#15709)

Created: 10/29/2018 02:38 PM (PT)

Shared: 01/15/2019 11:37 AM (PT)

---

This pre-registration is not yet public. This anonymized copy (without author names) was created by the author(s) to use during peer-review. A non-anonymized version (containing author names) will become publicly available only if an author makes it public. Until that happens the contents of this pre-registration are confidential.

---

### 1) Have any data been collected for this study already?

No, no data have been collected for this study yet.

### 2) What's the main question being asked or hypothesis being tested in this study?

- 1) If an uncertain outcome is associated with painful consequences, does that lead people to engage in wishful thinking by underestimating the probability of that outcome?
- 2) Can higher incentives for accuracy reduce wishful thinking?

### 3) Describe the key dependent variable(s) specifying how they will be measured.

There are two key dependent variables:

- 1) Accuracy: this is a binary variable (1:correct; 0:incorrect), indexing at each trial whether the participant correctly identified the displayed pattern.
- 2) Confidence: this is elicited trial-by-trial, as a number between 50% and 100%: after each choice participants are asked to report the probability that this answer is correct (50% is chance level, and 100% is certainty). On the basis of this we can construct a "belief" measure, which measures on a scale from 0 to 100% the belief of the subject about the orientation of the true pattern.

### 4) How many and which conditions will participants be assigned to?

In total, a participant sees 216 Gabor patches, and has to recognize whether these patterns are tilted to the right or left. After deciding between these two answers, the participant indicates his confidence in this decision in percentages. The participants are incentivized monetarily for providing accurate answers by a matching probability. At the start of the experiment, each participant is connected to an electric stimulation device that is personally calibrated to deliver mild but unpleasant shocks. Participants will receive an electric shock with a probability of 1/3 if the true answer is either right or left (depending on the condition).

This is a 2x2 design with 4 conditions. Each participant will participate in each condition (a within-subject design), and the pattern recognition tasks are equally divided over all conditions (54 in each condition). The two treatments dimensions are a) the incentives for accuracy (high and low), and b) whether the shock is associated with the left-leaning Gabor patch or the right-leaning Gabor patch.

### 5) Specify exactly which analyses you will conduct to examine the main question/hypothesis.

We are interested how the shock influences the accuracy and confidence of the participants. Specifically, we will test the following directional hypotheses:

- 1) Wishful thinking 1: Does the accuracy of identifying a given pattern go down if the potential shock is aligned with the true answer (the unpleasant true answer), relative to the case where the potential shock is not aligned with the answer? We use one-sided t-tests to evaluate the differences in the average accuracy and confidence between conditions, where each observation is the average of a subject's answers in that condition.
- 2) Wishful thinking 2: Similarly, we will use one-sided t-tests to examine whether the confidence in the true answer decreases if the potential shock is aligned with the true answer.
- 3) Accuracy incentives: We will test use one-sided t-tests whether accuracy and confidence in the true answer are higher in the condition with high incentives for accuracy.

In addition to t-tests, we will also use multivariate linear regression analysis to test the effect of our treatments (accuracy incentives, shock alignment) as well as their interaction. Finally, we will use linear mixed effect models (with or without individual fixed effects) where we can control for trial characteristics and/or subject characteristics (see below).

### 6) Describe exactly how outliers will be defined and handled, and your precise rule(s) for excluding observations.

We calibrate the difficulty of the task in the beginning of the experiment, so that we expect participants to be accurate 75% of the time on average. The actual accuracy in the experiment may deviate from this, and we will exclude the participant if actual accuracy is outside the [60%-90%] range, as this may indicate that, despite the calibration, the task was either too easy or too hard to detect meaningful differences.

### 7) How many observations will be collected or what will determine sample size? No need to justify decision, but be precise about exactly how the number will be determined.

Verify authenticity:<http://aspredicted.org/blind.php?x=mb5y37>

We ran a previous study that we had to discard due to an error in the code but was very similar. On the basis this study, running the same tests as specified above, we calculated that we could achieve more than 80% power with a sample of 60 people, so we will invite 60 participants

**8) Anything else you would like to pre-register? (e.g., secondary analyses, variables collected for exploratory purposes, unusual analyses planned?)**

We are likely to do some further exploratory analysis. For instance, we will see if the strength of the effect size differs by the difficulty of the task. We also look whether people who score higher on the trait anxiety, which we measure with a psychological questionnaire, are more affected by the shock, both in their accuracy and their beliefs. To investigate this, we will run mixed models which feature interactions between the shock treatment and the trait anxiety measured by the questionnaire.

Finally, we might explore how the effects of shocks play in typical models of confidence formation (inspired from signal-detection theory): confidence is known to be an increasing function of evidence for correct answers and decreasing for incorrect answers. We can test if the presence of shocks modulate the intercept of the slopes of this model (see e.g. Lebreton, et al. (2017) biorXiv for similar analysis)

## CONFIDENTIAL - FOR PEER-REVIEW ONLY

### Anticipatory anxiety about monetary losses and wishful thinking 2021 (#57718)

Created: 02/08/2021 03:07 AM (PT)

This is an anonymized copy (without author names) of the pre-registration. It was created by the author(s) to use during peer-review. A non-anonymized version (containing author names) should be made available by the authors when the work it supports is made public.

#### 1) Have any data been collected for this study already?

No, no data have been collected for this study yet.

#### 2) What's the main question being asked or hypothesis being tested in this study?

Our main question is about "wishful thinking": If an ambiguous state-of-the-world is associated with monetary losses, does that lead people to be less likely to correctly identify that state?

Secondary hypotheses: Does wishful thinking increase with...

- 1)  ...lower incentives for accuracy?
- 2)  ...higher monetary losses?
- 3)  ...increased ambiguity of the evidence?

#### 3) Describe the key dependent variable(s) specifying how they will be measured.

A participant sees up to 96 Gabor patches, and has to recognize whether these patterns are tilted to the right or left. The main dependent variable is based on their accuracy on this task in each trial, i.e. whether the participant correctly identified the displayed pattern. Wishful thinking is constructed as the difference in accuracy in recognizing patterns associated with monetary losses and those not associated with such losses.

#### 4) How many and which conditions will participants be assigned to?

Participants are endowed with an amount of money. They can lose part of this endowment on each trial with a probability of 1/3 if the true state is either right or left (depending on the condition) – i.e. monetary losses are not associated with participants' answers, but with the Gabor actual tilt direction. Additionally, participants are rewarded monetarily for providing accurate answers on the perception task.

There are 4 treatment dimensions. Each participant will participate in each condition (a within-subject design), and the pattern recognition tasks are equally divided over all treatments. The treatments dimensions are a) monetary loss is associated with the left-leaning Gabor patch or the right-leaning Gabor patch, b) the size of those losses (zero, low, or high), (c) the incentives for accuracy (high and low), and d) the ambiguity of the pattern (easy versus hard to discriminate).

#### 5) Specify exactly which analyses you will conduct to examine the main question/hypothesis.

We are interested how the losses influences the accuracy and confidence of the participants. Specifically, we will test the following directional hypotheses:

- 1) Wishful thinking (main hypothesis): Does the accuracy of identifying a given pattern go down if the potential loss is aligned with the Gabor patch reflecting the correct answer, relative to the case where the potential shock is not aligned with the answer? We use a t-test to evaluate the differences in the average accuracy and confidence between conditions, where each observation is the average of a subject's answers in that condition.
- 2) Secondary hypotheses: We will test use t-tests to assess whether wishful thinking is higher in the conditions with (a) low incentives for accuracy, (b) higher ambiguity of the pattern, and (c) higher potential losses.

In addition to t-tests, we will also use multivariate linear regression analysis to test the effect of our treatments (accuracy incentives, loss alignment) as well as their interaction. We will use linear mixed effect models where we can control for trial characteristics and/or subject characteristics. Finally, we will study the effect of incentives for accuracy on accuracy in the task.

#### 6) Describe exactly how outliers will be defined and handled, and your precise rule(s) for excluding observations.

Our experiment takes place on prolific. We will exclude participants who fail to answer simple attention checks at the beginning and throughout the experiment.

#### 7) How many observations will be collected or what will determine sample size? No need to justify decision, but be precise about exactly how the number will be determined.

We will recruit 220 subjects to the experiment.

#### 8) Anything else you would like to pre-register? (e.g., secondary analyses, variables collected for exploratory purposes, unusual analyses planned?)

More exploratory analysis will correlate wishful thinking with questionnaire items related to the self-reported anxiety of subjects about losing money, as

well as other questionnaire responses that include trait anxiety and emotion regulation.



## CONFIDENTIAL - FOR PEER-REVIEW ONLY

### Anticipatory anxiety about monetary losses, wishful thinking, and the effect (#83830)

Created: 12/21/2021 10:06 AM (PT)

This is an anonymized copy (without author names) of the pre-registration. It was created by the author(s) to use during peer-review. A non-anonymized version (containing author names) should be made available by the authors when the work it supports is made public.

#### 1) Have any data been collected for this study already?

No, no data have been collected for this study yet.

#### 2) What's the main question being asked or hypothesis being tested in this study?

In previous experiments we have found evidence for "wishful thinking". If a state of the world is associated with aversive outcomes, then subjects are less likely to correctly identify that state. Here, we will test whether wishful thinking is motivated by the threat of monetary losses. However, the main directional hypothesis we test is whether wishful thinking decreases with higher incentives for accuracy.

#### 3) Describe the key dependent variable(s) specifying how they will be measured.

A participant sees several short sequences of Gabor patches, and is asked to recognize whether these sequences are more right- or more left-tilted. The main dependent variable is based on their accuracy on this task in each trial, i.e., whether the participant correctly identified the displayed patterns. Wishful thinking is constructed as the difference in accuracy in recognizing patterns not associated with monetary losses and patterns associated with such losses. Our paradigm expands on an earlier experiment, which focused on the accurate recognition of the tilt of a single Gabor patch.

#### 4) How many and which conditions will participants be assigned to?

Participants are endowed with an amount of money. They can lose part of this endowment on each trial with a probability of 1/3 if the true state is either right or left tilted (depending on the condition) – i.e., these monetary losses are not associated with participants' answers, but with the actual tilt of the sequence. In addition, participants are rewarded with a monetary bonus for correctly identifying sequences.

There are 3 treatment dimensions. Each participant will participate in each condition (a within-subject design) and the pattern recognition tasks are equally divided across treatments. The treatment dimensions are a) monetary loss is associated with a left-leaning sequence of Gabor patches or the right-leaning sequence, b) the incentives for accuracy are high or low, and c) the ambiguity of the sequences of patterns is high or low. Our measure for ambiguity derives from a continuous variable of pattern difficulty that we dichotomize. The experimental uses a 2 x 2 x 2 within-subjects design with the factors loss pattern (aligned, not aligned with loss), incentives (high, low) and ambiguity (high, low).

#### 5) Specify exactly which analyses you will conduct to examine the main question/hypothesis.

We will run a linear regression of correct answers on the following variables

- [loss pattern]: a dummy for whether the pattern is associated with a loss.
- [high incentives]: a dummy for the high incentive condition.
- [interaction]: a dummy for the interaction of the two conditions.

We cluster standard errors at the individual level. We run regressions where each observation is the individual average over trials in that condition, as well as regressions where we use all data points.

We test for

- wishful thinking by testing whether the first coefficient is positive and stat. significant.
- the effect of incentives on overall task performance by testing whether the second coefficient is positive and stat. significant.
- the interaction (our main hypothesis), by testing whether the third coefficient is positive and stat. significant.

For patterns associated with a loss, incentives should raise performance both because of the overall effect and because of a reduction in wishful thinking. To test this joint effect, we conduct an additional t-test assessing whether high accuracy incentives improve performance on loss patterns only.

#### 6) Describe exactly how outliers will be defined and handled, and your precise rule(s) for excluding observations.

We conduct our online experiment on Prolific. We will exclude participants who fail to answer simple attention checks at the beginning and throughout the experiment.

#### 7) How many observations will be collected or what will determine sample size? No need to justify decision, but be precise about exactly how the number will be determined.

400

#### 8) Anything else you would like to pre-register? (e.g., secondary analyses, variables collected for exploratory purposes, unusual analyses planned?)

We will test whether self-reported concentration responds to incentives, as a manipulation check.

We will test our secondary hypothesis that wishful thinking increases with the ambiguity of the evidence.

Leveraging our within subject identification of wishful thinking, we will also study the heterogeneity of wishful thinking across subjects, correlating it with responses from the exit questionnaire.

We will test the robustness of our results if we exclude participants whose average accuracy is worse than chance.

## CONFIDENTIAL - FOR PEER-REVIEW ONLY

### Anticipatory anxiety, wishful thinking, and the effect of accuracy incentives in a dot-task. (#89876)

Created: 03/04/2022 02:58 AM (PT)

This is an anonymized copy (without author names) of the pre-registration. It was created by the author(s) to use during peer-review. A non-anonymized version (containing author names) should be made available by the authors when the work it supports is made public.

#### 1) Have any data been collected for this study already?

No, no data have been collected for this study yet.

#### 2) What's the main question being asked or hypothesis being tested in this study?

In previous experiments we have found evidence for "wishful thinking". If a state of the world is associated with aversive outcomes, then subjects are less likely to correctly identify that state.

Here, we will test whether wishful thinking is motivated by the threat of monetary losses in a task we have not previously used. However, the main directional hypothesis we test is whether wishful thinking decreases with higher incentives for accuracy.

#### 3) Describe the key dependent variable(s) specifying how they will be measured.

On each trial, a participant sees an array of 100 blue and red dots, and is asked to recognize whether there are more blue or red dots. The main dependent variable is based on their accuracy on this task in each trial, i.e., whether the participant correctly identified the displayed patterns. Wishful thinking is constructed as the difference in accuracy in recognizing patterns not associated with monetary losses and patterns associated with such losses.

Our paradigm expands on an earlier experiment, which focused on the accurate recognition of the tilt of Gabor patches.

#### 4) How many and which conditions will participants be assigned to?

Participants are endowed with an amount of money. They can lose part of this endowment on each trial with a probability of 1/3 if the true state is either right or left tilted (depending on the condition) – i.e., these monetary losses are not associated with participants' answers, but with the actual tilt of the sequence. In addition, participants are rewarded with a monetary bonus for correctly identifying sequences.

There are 3 treatment dimensions. Each participant will participate in each condition (a within-subject design) and the pattern recognition tasks are equally divided across treatments. The treatment dimensions are a) monetary loss is absent or present, and when it is present, it associated with patterns with a majority of red dots, or with a majority of blue dots, b) the incentives for accuracy are high or low, and c) the asymmetry in the number of blue and red dots differs. Our measure for ambiguity derives from a continuous variable of pattern difficulty that we dichotomize. The experiment uses a 3 x 2 x 4 within-subjects design with the factors loss pattern (no losses, associated with blue patterns, associated with red patterns), incentives (high, low) and ambiguity (high, medium high, medium low, low).

#### 5) Specify exactly which analyses you will conduct to examine the main question/hypothesis.

We will run a linear regression of correct answers on the following variables

1. [loss pattern]: a dummy for whether the pattern is associated with a loss. A t-test on the coefficient tests our hypothesis on wishful thinking.
2. [high incentives1]: a dummy for the high incentive condition when there is no possibility of monetary losses. A t-test on the coefficient tests for the effect of high incentives in the absence of losses.
3. [high incentives2]: a dummy for the high incentive condition when there is a possibility of monetary losses. A t-test on the coefficient tests for the effect of high incentives when losses are present.
4. [interaction]: restricting ourselves to the pattern with losses, we will look at a dummy for the interaction of the loss patterns with the high incentive condition. A t-test on the coefficient for the interaction term tests whether the incentives reduce wishful thinking.

We cluster standard errors at the individual level. We run regressions where each observation is the individual average over trials in that condition, as well as regressions where we use all data points.

For patterns associated with a loss, incentives should raise performance both because of the overall effect and because of a reduction in wishful thinking. To test this joint effect, we conduct an additional t-test assessing whether high accuracy incentives improve performance on loss patterns only.

#### 6) Describe exactly how outliers will be defined and handled, and your precise rule(s) for excluding observations.

We conduct our online experiment on Prolific. We will exclude participants who fail to answer simple attention checks at the beginning and throughout the experiment.

#### 7) How many observations will be collected or what will determine sample size? No need to justify decision, but be precise about exactly how the number will be determined.

We will recruit 400 subjects to complete the experiment.

**8) Anything else you would like to pre-register? (e.g., secondary analyses, variables collected for exploratory purposes, unusual analyses planned?)**

We will test whether self-reported concentration responds to incentives, as a manipulation check.

We will test whether self-reported anxiety responds to the presence of losses, as a manipulation check.

We will test our secondary hypothesis that wishful thinking increases with the ambiguity of the evidence.

Leveraging our within subject identification of wishful thinking, we will also study the heterogeneity of wishful thinking across subjects, correlating it with responses from the exit questionnaire.

We will test the robustness of our results if we exclude participants whose average accuracy is worse than chance.

## **D Exclusion based on instruction questions**

All experiments included quiz questions to check participant understanding of the instructions (see E).

In Experiments 2-4, the quiz questions are presented intermixed with the instructions, separated in two sections. Each question is repeated upon a wrong answer, and each section is repeated up to 2 times if 3 or more mistakes are made within that section. In Experiment 2 the first section was repeated up to 3 times, though this rarely happened. Participants were excluded from the experiment if they made more than 4 mistakes in total.

## **E Experimental Instructions**

Instructions for the different Experiments can be found on the Online Supplementary Material: <https://osf.io/tznpy/>.

## **F Attention Checks**

Several attention checks (see screen captures on the Online Supplementary Material: <https://osf.io/gn9uk>) were used at the beginning of the task for Experiments 2-4. Each was repeated upon a wrong response, and the experiment stopped at two total wrong responses. The last attention check (pressing a key written on screen) appeared once again at the end of the first block, and the experiment ended if answered incorrectly four times.