

# Minding the Gap:

## On the Origins of the Description-Experience Gap\*

Ryan Oprea<sup>†</sup> Ferdinand M. Vieider<sup>‡</sup>

December 5, 2025

### Abstract

We provide evidence that “noisy coding” is responsible for both (i) classic probability-dependence of risk-taking and (ii) its reversal when the properties of lotteries are learned by sampling rather than by explicit description. Guided by a stylized model of noisy sampling, we show that simply forcing experimental subjects to sample redundant information about the primitives of lotteries causes both types of probability-dependence to disappear, closing the description-experience gap and resulting in broadly neoclassical behavior. This strongly suggests that these anomalies are a joint outgrowth of decision makers’ noisy representations of the primitives of lotteries rather than expressions of true risk preferences.

**Keywords:** risk taking; noisy coding; probability weighting; decision from experience

**JEL codes:** C91, D91, G0

## 1 Introduction

Risk-taking has been extensively documented to deviate from the predictions of the standard model of expected utility theory (*EUT*). A key anomaly identified in the last half century

---

\*We are grateful to Larbi Alaoui, Rava Azeredo da Silveira, Ido Erev, Cary Frydman, Richard Gonzalez, Ralph Hertwig, Alex Imas, Sebastian Olschewski, to participants at the Economics Psychology Seminar in Basel, the EVRE Workshop in Grenoble, the 2023 RISL $\alpha\beta$  Workshop, the Workshop in Honor of Peter Wakker, and the Summer School on the “Cognitive Foundations of Decision-Making” for helpful comments and suggestions. All errors remain our own. This research was supported by the National Science Foundation under Grant SES-1949366 and by the Research Foundation Flanders under the project “Causal Determinants of Preferences” (G008021N). It was approved by UC Santa Barbara IRB.

<sup>†</sup>Oprea: Economics Department, University of California, Santa Barbara, roprea@gmail.com.

<sup>‡</sup>Vieider: RISL $\alpha\beta$ , Economics Department, Ghent University, Belgium, fvieider@gmail.com.

consists in systematic *probability-dependence* of risk-taking. As hundreds of experiments have shown, experimental subjects, when given explicit descriptions of lotteries, tend to be more risk-taking for small probabilities and relatively less risk-taking for large probabilities of winning a prize. This probability-dependence has been enshrined as a centerpiece of alternatives to EUT such as prospect theory (Kahneman & Tversky 1979, Tversky & Kahneman 1992). When subjects are required to discover the properties of lotteries by sampling from them instead of by reading explicit descriptions of their properties, the direction of this anomaly reverses: subjects now tend to be highly risk averse for low probabilities, and risk-taking tends to *increase* in the probability of winning a prize (Barron & Erev 2003, Hertwig et al. 2004). This constitutes an important paradox under prevalent EUT-alternatives such as prospect theory. To date, these two patterns lack a unified theoretical explanation.

Understanding the reversal in the dependence of risk-taking propensities on probabilities when choice options are described or need to be experienced thus appears as a key ingredient for understanding what drives risk-taking in general. Experience-based choice has been largely explained through sampling error (Hertwig et al. 2004, Fox & Hadar 2006, Hertwig & Pleskac 2010), but removing sampling error has failed to eliminate the gap (Hau et al. 2008, Ungemach et al. 2009). Here, we propose a unified explanation of opposite probability-dependence of risk-taking in described versus experience-based choices based on “noisy cognition”. Noisy cognition has recently been proposed as an explanation of standard probability-dependence in described choices (Zhang et al. 2020, Enke & Graeber 2023, Oprea 2024, Vieider 2024b, Frydman & Jin 2025, Khaw et al. 2025). We extend this theoretical framework to a sampling-based setting, thereby showing that 1) noisy cognition provides a unified theoretical setup under which to rationalize description- and experience-based choices; 2) by leveraging the insights from the model on the causes of probability-dependence in DfD, we are able to experimentally remove such standard probability-dependence; and 3) by combining this treatment with a similar intervention on DfE, we can finally close the the gap between description-based choice and experience-based choice.

**The description-experience gap and its significance.** Suppose a decision maker (DM) has to make a choice between a sure amount  $c$  and a lottery that pays  $x > c$  with probability  $p$  (and  $y < c$  otherwise).<sup>1</sup> In what has come to be the standard protocol, DMs are explicitly told how many outcomes each lottery can produce, the payoffs each outcome results in and

---

<sup>1</sup>We will use this simple choice as a running example, and our experiment will exclusively employ such simple choices. Our framework extends to losses in a straightforward manner. It can also be extended to multi-outcome lotteries via an N-dimensional generalization; see the Online Appendix for details.

the probabilities of each outcome. The DM uses this information to choose the lottery she prefers. Call this standard paradigm “decision from description”, or *DfD*. More recently, researchers have studied an alternative paradigm to DfD for studying lottery choice. In “decisions from experience” (DfE) experiments (Barron & Erev 2003, Hertwig et al. 2004), subjects are told nothing about the two lotteries but must learn all of their properties entirely by *sampling* each of them. In standard DfE experiments (under the so-called “sampling paradigm”), subjects choose how many times to sample each lottery and use the information gleaned from these samples to make their decision.

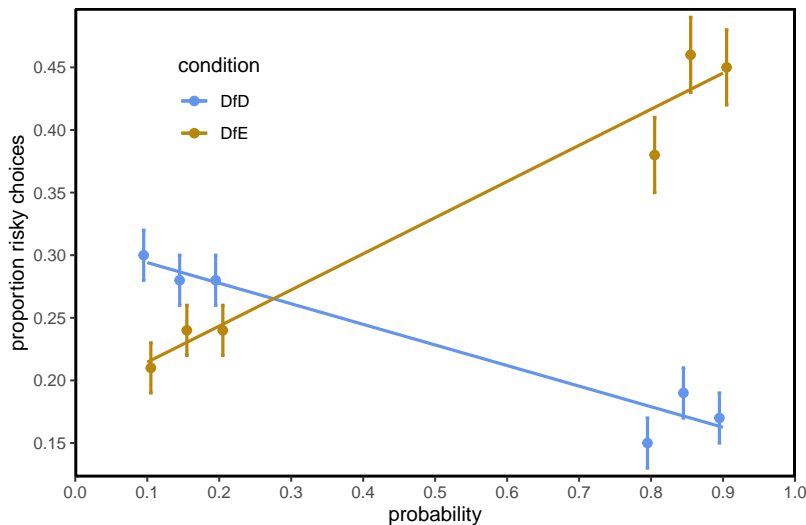


Figure 1: The GAP: Decisions from Description vs Decisions from Experience

The figure shows choice proportions for the risky lottery in DfD and DfE, ordered by probability of winning. The figure shows between subject comparison based on identical tasks that are either described (DfD) or sampled (DfE), fitted with a linear regression line. The tasks are constructed in such a way that the sure amount varies symmetrically around the expected value of the lottery (see below for details). The error bars indicate  $\pm 1$  standard error.

Figure 1 illustrates the opposite changes in risk-taking over the probability interval one observes in DfD vs DfE based on a replication experiment we conducted (see below for details). When facing lotteries with a small probability of winning, DMs take more risk in DfD than in DfE. This tendency, however, completely flips for large probabilities of winning: subjects are now much more risk-taking in experience-based choice. The inverted responses in DfE and DfD produce what the literature has called the “decision-experience gap” (hereafter, simply the *GAP*) in lottery choice. The significance of this GAP stems from the observation that it remains an open mystery in the literature. In particular, it has eschewed uniform modeling — constituting a paradox under behavioral models of risk-taking such as prospect theory — and has resisted repeated attempts at closing it

(reviewed below). Understanding the source of this gap is crucial to understanding not just probability-dependence, but the determinants of risk-taking more generally.

**Our contribution.** In this paper we offer a theoretical explanation for the GAP that also explains the nature of probability-dependence in risk-taking. Our explanation is rooted in a kind of irony: the decision-experience gap, we argue, is a consequence not only of the fact that DfD and DfE are psychologically different, but also of the fact that they are in an important sense more psychologically similar than has been previously recognized. Drawing on arguments and evidence from neuroscience, we argue that the kind of *explicit* sampling that occurs in DfE also necessarily occurs *implicitly* in the brain when a subject reasons about the properties of fully described lotteries in DfD. The noisy representations arising from finite sampling can simultaneously generate classic probability-dependence as observed in DfD and its inversion in DfE.

The inversion in DfE, in particular, is driven by the role of sampling variance in endogenously determining when to stop sampling. This, in turn, will drive the extent of sampling *error* — a key element for understanding experience-based choices (Hertwig et al. 2004, Fox & Hadar 2006, Hertwig & Pleskac 2010). A key novelty we present here is the endogenous nature of the decision on when to stop sampling, which allows us to provide a unified representation of decision processes at work in DfD and DfE. The prediction arising from this is striking: ex ante risk averse DMs should sample more from lotteries with large probabilities than from those with small probabilities. This will reduce sampling error for large probabilities while concomitantly increasing confidence in the sampled proportions, thus inducing ex ante risk averse DMs to take *more risk* for large probabilities in DfE. This mechanism is indeed strongly supported in our data.

Our theoretical insights also allow us to explain why previous efforts to eliminate the GAP have failed. In particular, we show that efforts to close the GAP by forcing subjects in DfE to sample more intensively than they naturally would — as done e.g. by Hau et al. (2008, 2010), Ungemach et al. (2009), Aydogan & Gao (2020), Cubitt et al. (2022) — has the unexpected effect of simultaneously removing the sampling variance needed for standard probability-dependence to occur, leading to behavior broadly consistent with EUT, rather than standard probability-dependence.<sup>2</sup> This means the GAP can never be eliminated by

---

<sup>2</sup>Note that our results here are fully consistent with previous attempts at closing the GAP in similar ways, such as the one of Ungemach et al. (2009). Just like shown by the latter, forced sampling in DfE alone does not close the GAP. By using richer choice tasks, however, we can directly test for probability-dependence in risk-taking, something that previous studies could not do directly.

forcing subjects in DfE to observe larger samples alone: in order to close the GAP, we must also increase the precision of neurally coded probabilities in DfD, and thus reduce classical probability-dependence in described choices.

Our contribution here is to offer a particularly direct type of evidence for this hypothesis, and to show that it accounts for the description-experience gap. The main novelty on the experimental side thus emerges from a treatment that requires subjects to take large, balanced samples from fully described choice options. Even though such samples are completely redundant from the perspective of preference-based models such as prospect theory, they allow us to address the fundamental cause underlying noisy cognition. Our results indicate that such samples have a large effect: after forced samples subjects make mildly risk averse lottery choices that broadly comply with standard EUT. By applying a similar manipulation to DfE we further show that eliminating the causes predicted by the model to underlie the opposite deviations from EUT allows us — for the first time ever — to completely close the description-experience GAP.

We finally supplement this evidence with an additional treatment that allows — but does not force — subjects to freely sample from described choices. This allows us to show, first of all, that subjects do indeed feel a need for sampling, even for fully described choices and in the presence of considerable opportunity costs. Just as importantly, our model predicts this treatment to introduce sampling error into description-based choice. Our data show that this is clearly the case. By letting subjects freely sample from fully described choice options, we introduce positive probability-dependence of the DfE kind into DfD. Although some differences in the degree of probability-dependence remain, this achieves something that acting on DfE alone has never achieved: we can close the GAP by acting on one of the two experimental paradigms alone.

**Fit with the literature.** Our paper contributes to several literatures. First is a long running literature on probability-dependence in risk-taking and related anomalies, going back to Preston & Baratta (1948). Such probability-dependence became a key component of prospect theory, where it is captured by an inverse S-shaped probability weighting function (Kahneman & Tversky 1979, Tversky & Kahneman 1992, Tversky & Wakker 1995, Wakker 2010), and is the mechanism by which that theory accounts for phenomena like the coexistence of lottery play and insurance uptake and the Allais paradoxes. Numerous empirical studies have documented systematic increases of relative risk aversion in the probability of winning a prize — Imai et al. (2025) provide a meta-analytic overview of the evidence in

DfD.

The second is a literature documenting the gap between DfD and DfE (Barron & Erev 2003, Hertwig et al. 2004). Sampling error was proposed as an early explanation for the GAP (Fox & Hadar 2006). However, subsequent investigations showed that, although sampling error is an important contributor to the GAP, interventions including (i) eliminating sampling error by matching probabilities in DfD to DfE, (ii) increasing the samples by offering higher stakes, and (iii) forcing people to sample the complete urn in DfE fail to eliminate the gap (Hau et al. 2008, Ungemach et al. 2009, Hau et al. 2010, Hertwig & Pleskac 2010, Wulff et al. 2018). Because of this, the underlying causes of the GAP have largely remained a mystery — see Hertwig & Erev (2009) and de Palma et al. (2014) for narrative reviews, and Wulff et al. (2018) for a systematic meta-analysis of the decision-experience gap and possible factors contributing to it. Cubitt et al. (2022) present a careful experimental decomposition, which concludes that sampling error is the prime driver of the GAP (but once more fails to close the GAP by eliminating sampling error, pointing at missing pieces in the explanation).

The third is a growing literature documenting the role noisy cognition plays in behavioral anomalies (Natenzon 2019, Khaw et al. 2021, Frydman & Jin 2022). Most closely related is a line of research examining how cognitive noise (and efficient ways the brain deals with such noise) contributes to distorted perceptions of probabilities (Zhang & Maloney 2012, Steiner & Stewart 2016, Zhang et al. 2020, Enke & Graeber 2023, Herold & Netzer 2023, Netzer et al. 2024, Frydman & Jin 2025, Khaw et al. 2025, Oprea 2024, Vieider 2024*b*). More broadly, our work is related to a literature documenting the role cognitive frictions play in decision-making under risk (Enke & Graeber 2023, Bohren et al. 2024, Oprea 2024) and, broader still, the way cognitive constraints and the brain’s response to these constraints explain a wide class of anomalies in decision-making (Simon 1959, Robson 2001*a,b*, Netzer 2009, Robson & Samuelson 2011).<sup>3</sup>

---

<sup>3</sup>A recent, contemporaneous paper, Bohren et al. (2024), documents and decomposes a complementary description-experience gap that operates in richer environments than the one we (and the previous description-experience literature) study. In evaluating realistic lotteries with many potential outcomes (e.g., eleven states), they show that subjects’ behavior tends to be constrained by memory limitations in DfE, while it tends to be constrained by attentional limitations in DfD. This leads to systematic differences in lottery choices in DfE and DfD environments – a gap that can be eliminated with aids to attention and memory. Memory seems to play much less of a role when studying simpler choice situations such as the ones we use here — see the GAP decomposition by Cubitt et al. (2022) for details.

## 2 Drivers of risk-taking in DfE

What is responsible for the GAP and the reversals of probability-dependence that produce it? Here we start from DfE, highlighting two basic features of the information structure. The first is the well-known *sampling error*: unless the DM collects a very large sample, she runs the risk of drawing misleading samples that systematically distort beliefs particularly at extreme probabilities. The second (which has not been emphasized in the literature so far) we will call *sampling variance*: because the DM’s sample is finite, she cannot be entirely confident in the sample she draws. This will make it optimal to combine such samples with her prior beliefs in a Bayesian fashion, distorting her posterior beliefs. As we will show, these two features interact, and are responsible for the positive probability-dependence of risk-taking observed in DfE.

To fully specify a model of DfE, we must describe not only how people form beliefs about probabilities and payoffs, but also how these beliefs co-evolve with higher order beliefs about the structure of the lotteries (e.g., the number of outcomes in each lottery’s support). To close the model, it is therefore necessary to make a number of detailed modeling choices about the evolution of these structural beliefs that do not directly impact the way we interpret and design our experiments. In the Online Appendix A.1 we propose such a fully specified model.<sup>4</sup> But in this section, for expositional ease, we abstract from these issues of higher order belief formation altogether by (i) assuming that subjects already know the structure of the lotteries<sup>5</sup>, (ii) that subjects quickly identify which lottery is risky during sampling and (iii) by focusing attention on the way subjects evaluate the risky arm. In the fully specified, general model in Online Appendix A we discuss the implications of these assumptions, but argue that they are qualitatively irrelevant to the key matters at hand.

**Basic model structure.** To model the way beliefs change as a DM samples the simple binary lotteries in our experiment, let  $\alpha$  be the number of draws in which the DM observed payment  $x$  and  $\beta$  the number of draws in which she observed payment  $y$ . We model

---

<sup>4</sup>In the full version of the model in Online Appendix A, we close the model by assuming that (i) subjects mainly use samples to build beliefs about the comparative properties of the two choice options (which seems likely given the choice subjects face), (ii) that subjects know that they are making a risky choice and that the choice is therefore not between two degenerate lotteries (which seems likely given the lotteries subjects exclusively see in Part 1 of the experiment) and (iii) make a few other technical assumptions required to fully specify the joint inference problem. The main implication of (ii) is that inferences in which the outcomes observed in both choice options are attributed probability close to 1 will carry very high noise, in a sense to be made precise below. Within the formalism of the model, this assumption mainly serves to explain why subjects take more than 1 sample from each option.

<sup>5</sup>In our experiment, this is in fact a fairly realistic assumption, given that subjects entering DfE have all just made a number of lottery choices, all with the same structure.

this sampling process using a Beta distribution with parameters  $\alpha$  and  $\beta$ , producing a representation of the probability  $p$  of earning  $x$  equal to  $\mathbb{E}[\hat{p}|p] = \frac{\alpha}{\alpha+\beta}$  (i.e. the sampled mean probability  $\hat{p}$ , given the true probability  $p$ ).<sup>6</sup> We will assume that the DM’s beliefs are represented in a log-odds form. This is not necessary for any of our qualitative conclusions in what follows, but (i) it is increasingly supported in neuroscience both empirically and theoretically<sup>7</sup> and (ii) it will allow us to neatly connect our characterization to a linear in log-odds (LLO) functional form that is commonly used to characterize probability-dependence in risk-taking in the prospect theory literature (Gonzalez & Wu 1999).

**Sampling error.** Successes  $\alpha$  and failures  $\beta$  are, on average, sampled in an unbiased way (i.e.,  $\ln\left(\frac{\alpha}{\beta}\right) = \ln\left(\frac{p}{1-p}\right)$  on average). However, the binomial distribution will produce samples that underestimate small probabilities and overestimate large probabilities (unless the samples are unrealistically large). This issue — typically termed *sampling error* in the DfE literature — has been discussed as one of the key drivers of the GAP from the very beginning (Hertwig et al. 2004). Fox & Hadar (2006) argued that sampling error may indeed be the *sole* driver of the GAP, and that eliminating it ought to result in choice patterns that converge towards those observed in DfD. The subsequent literature has thus devoted much energy to trying to eliminate sampling error, either by incentivizing or forcing subjects to take larger samples (Hau et al. 2008), or by forcing subjects to take large, balanced samples from both choice options (Ungemach et al. 2009), or by only selecting samples that happen to reflect the true underlying probability (Wulff et al. 2018). Two key insights resulting from this literature are that 1) sampling error explains at least part of the GAP; but 2) while reducing or eliminating sampling error narrows the GAP, it fails to close it completely.

---

<sup>6</sup>It is important to emphasize that our model *does not require us to assume* that the DM knows the structure of the decision problem. We use a Beta distribution here purely for expositional simplicity, and because binary lotteries is all a DM will ever experience in our experiments. Our model generalizes to any number of outcomes by using a Dirichlet distribution—the multi-dimensional generalization of the Beta—to represent the different states. Indeed, we can use Dirichlet distributions defined over all possible outcomes to explicitly model the inference process of the DM about the underlying state space in DfE—an important element that distinguishes our approach from some of the DfE literature in economics, which has assumed that the DM (often counterfactually) knows the objective state space or which has (in some papers) provided this information ex ante in experiments (Abdellaoui et al. 2011, Aydogan 2021, Cubitt et al. 2022). Online Appendix A provides details of the inference process, and of how the model we use here can be generalized to  $N$  states of nature.

<sup>7</sup>It is common in neuroscience to assume that the brain represents the sort of evidence encoded by  $\alpha$  and  $\beta$  in terms of log-odds. This is in part because of its computational efficiency for the brain, a straightforward consequence of the fact that new evidence can be simply added to pre-existing evidence, which is a much less computationally expensive operation than, e.g., multiplication. It is also in part because of the empirical success of such representations. For instance, Zhang & Maloney (2012) describe log-odds representations as “ubiquitous”, discussing a long list of findings which can be fit by log-odds representations. Glanzer et al. (2019) identify a unique empirical signature of log-odds representations, and argue that such representations underlie neural representations in general.



Unbiased samples in any given task will only obtain if the DM takes very large (technically: infinite) samples. As a result, we should expect the ratio of  $\alpha$  and  $\beta$  observed by subjects in finite samples to produce systematically distorted impressions of the log odds. This is particularly true of samples taken from lotteries with extreme probabilities, where sampling error is most likely and where the gap between description and experience is most severe.

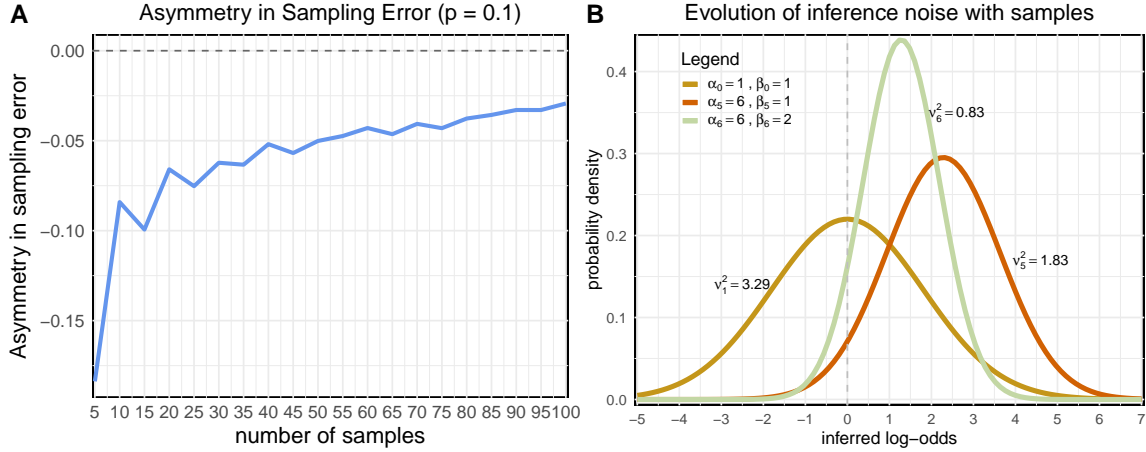


Figure 2: Sampling error and inference error in DfE

The figures show how sampling error and sampling variance evolve as samples accumulate. Panel A shows the likelihood to underestimate a probability of  $p = 0.1$  as the number of samples increases (by steps of 5 samples). The smallest number of samples for which such a probability can be accurately estimated is 10. At 10 samples, however, the likelihood of underestimating the probability still exceeds the likelihood of estimating it correctly or larger by 7 percentage points. This asymmetry is reduced at a decreasing rate as samples increase. Panel B illustrates the evolution of sampling variance as a function of samples taken. Subsequent samples do not only update the “best guess” of the probability, but also reduce sampling variance.

Panel A in figure 2 illustrates the sampling error occurring in small samples (see also Hertwig & Pleskac 2010, for an extensive discussion). The figure illustrates the error asymmetry — the excess likelihood of sampling a probability that is smaller than the true probability compared to a probability that is larger — for  $p = 0.1$ . At 10 samples (the first number that can theoretically result in a correct estimate), a DM is still 8.7 percentage points ( $pp$ ) more likely to underestimate the true probability than to overestimate it. This error asymmetry subsequently decreases at a decreasing rate. At 100 samples (the number of non-representative samples imposed by Hau et al. 2008 in their experiment 3), the asymmetry in the direction of underestimation is still about 3pp. This highlights that 1) some degree of sampling error is almost inevitable for extreme probabilities in realistic samples; and 2) returns to sampling decrease rapidly once one exceeds a certain threshold.

**Sampling variance.** Because  $\alpha$  and  $\beta$  are finitely sampled, they produce noisy beliefs about the true probabilities. Note that such beliefs will be noisy in finite samples even if

the samples are accurate on average, and the DM can never be 100% sure of whether a *given* sample correctly reflects the underlying outcome-generating probability. This addition of sampling variance to the explanation of the GAP is a novel contribution we bring to the literature — and as we will see shortly — will provide the key to obtaining a unified understanding of the opposite types of probability-dependence observed in DfE and DfD.<sup>8</sup>

Given samples of successes  $\alpha$  and failures  $\beta$ , the “best guess” of the log-odds will be given by  $\ln\left(\frac{\alpha}{\beta}\right)$ , which is the log-odds equivalent of the mean of the Beta distribution (the sampled proportion of successes  $x$ ). The beliefs, however, must be augmented by an error term  $\varepsilon$  to capture variability in small samples. The log-odds formulation gives rise to approximately normally distributed errors even with relatively few observations (see e.g. Gelman et al. 2014, section 5.6), so that we will assume a logit-normal distribution for the errors. Following the characterization of the logit-normal distribution by Atchison & Shen (1980), it is straightforward to obtain an explicit solution for the sampling variance from the draws representing the odds:

$$\varepsilon \sim \mathcal{N}(0, \nu_n^2) \text{ , } \nu_n^2 = F'(\alpha_n) + F'(\beta_n), \quad (1)$$

where  $F'$  represents the trigamma function, and where we now subscript the samples and sampling variance by the number of samples  $n$  to emphasize the dependence of these quantities on the number of samples taken. The sampling *precision*  $\nu_n^{-2}$  (i.e. the inverse of the sampling variance) will increase in the number of draws  $n$ . We can thus interpret the precision  $\nu_n^{-2}$  as a measure of confidence in the sampled proportion  $\ln\left(\frac{\alpha_n}{\beta_n}\right)$ .

Panel B in figure 2 illustrates how sampling variance evolves with subsequent samples. Let us assume a DM starts sampling from initial parameters  $\alpha_0 = \beta_0 = 1$ . This corresponds to an ignorance prior with a uniform distribution attributing equal ex ante likelihood to all probabilities (i.e. Laplace’s rule of succession).<sup>9</sup> This distribution is centered on  $p = 0.5$ , but shows low confidence in that estimate (all probabilities are seen as equally likely). The distribution parameterized by  $\alpha_5 = 6$  and  $\beta_5 = 1$  shows the situation after 5 samples, all of which have yielded draws of the prize  $x$ . As one would expect, the mean estimate

---

<sup>8</sup>Olschewski & Scheibehenne (2024) present a discussion of different types of noise arising when DMs need to infer (and bet on) means of a series of sampled numbers, and present a concept of “Thurstonian uncertainty” that resembles what we here call sampling variance.

<sup>9</sup>Note that this specific value only serves illustrative purposes, and is in no way essential to our conclusions. This will become apparent shortly, when we will describe the Bayesian integration of the evidence from the sample with prior expectations. Online appendix A.1 further discusses the likelihood in a more general setting based on Dirichlet distributions used to infer the structure of the decision problem jointly with the probability attached to each state.

of the probability is now larger. Just as importantly, the distribution has narrowed — the sampling variance has decreased, thus increasing the ‘confidence’ the DM has in the sampled proportion. Assume now the DM draws a sample of  $y$ . This reduces the sampled proportion  $\alpha/\beta$ , but also further increases the precision of the sample. This illustrates an important property of the model: sampling variance is a decreasing function of the number of samples, but constitutes a conceptually separate dimension from the specific samples drawn.<sup>10</sup>

**Optimal Bayesian Inference.** Taking into account sampling variance in any given (finite) sample, a Bayesian DM will rationally combine the results of her sampling with her prior beliefs to draw inferences about the true underlying probability. Continuing with our log-odds characterization of beliefs, assume that the prior, too, takes logit-normal form:

$$\ln\left(\frac{p}{1-p}\right) \sim \mathcal{N}\left(\ln\left(\frac{p_0}{1-p_0}\right), \sigma^2\right). \quad (2)$$

As we show in more detail in Appendix A.2, the posterior expectation of the log-odds being inferred,  $\ln\left(\frac{\hat{p}}{1-\hat{p}}\right)$ , conditional on the true log-odds, will take the following form:

$$\begin{aligned} \mathbb{E}\left[\ln\left(\frac{\hat{p}}{1-\hat{p}}\right) \mid \ln\left(\frac{p}{1-p}\right)\right] &= \gamma_n \ln\left(\frac{\alpha_n}{\beta_n}\right) + (1-\gamma_n) \ln\left(\frac{p_0}{1-p_0}\right) \\ &= \ln\left(\frac{\alpha_n}{\beta_n}\right) + (1-\gamma_n) \left[\ln\left(\frac{p_0}{1-p_0}\right) - \ln\left(\frac{\alpha_n}{\beta_n}\right)\right], \end{aligned} \quad (3)$$

where  $\gamma_n = \frac{\sigma^2}{\sigma^2 + \nu_n^2}$  is the Bayesian evidence weight, i.e. the weight put on the sampled proportion relative to the prior expectation  $\ln\left(\frac{p_0}{1-p_0}\right)$ . Importantly,  $\gamma_n$  is itself an increasing function of the number of samples  $n$  taken, since it is inversely proportional to the sampling variance  $\nu_n^2$ . The equation above gives us some interesting intuitions (technical details in Online Appendix A). Defining  $\delta_n \triangleq \left(\frac{p_0}{1-p_0}\right)^{(1-\gamma_n)}$  and substituting it into the first line in (3) yields a linear in log-odds probability weighting function as often used in prospect theory (Gonzalez & Wu 1999).<sup>11</sup> The second line in (3) illustrates why this results in *biased inferences*: even for choice proportions  $\alpha_n/\beta_n$  that are on average correct, the regression to the mean of the prior will systematically distort the inferences drawn.<sup>12</sup>

<sup>10</sup>Technically this orthogonality is not perfect, since both will depend on the number of samples drawn to some extent. The dependence will be particularly strong for very small and extreme samples, e.g. when only one single outcome has been observed.

<sup>11</sup>This conclusion holds because on average the sampled proportions  $\alpha/\beta$  will be equal to the true odds.

<sup>12</sup>The expression indeed shows the definition of bias, inasmuch as it illustrates regression to the mean of the prior as a source of systematic deviations from  $\ln\left(\frac{\alpha_n}{\beta_n}\right)$ , which is an unbiased estimator of  $\ln\left(\frac{p}{1-p}\right)$  on average. It is important to note that — notwithstanding this systematic bias — the inference process is *optimal* given some constraints on sampling (e.g. in the presence of opportunity costs or time pressure, both

**The discriminability equation.** One key insight from the model, which will allow us to characterize behavior in DfE, derives from the observation that sampling variance and sampling error will interact. In particular, sampling variance — and the resulting regression to the mean of the prior described above — will determine at what point a DM concludes that she has sufficient information to stop sampling. This constitutes a key innovation of our approach, given that the DfE literature has paid relatively little attention to the decision when to stop sampling.<sup>13</sup> Here we will show that the sampling-stopping decision is endogenous to 1) the prior expectation of the DM; and 2) the precision of the sample drawn. The decision on when to stop sampling will in turn determine sampling error, and thus the positive probability-dependence or risk-taking observed in DfE.

To understand why DMs tend to undersample in DfE, we leverage the result on probabilistic inferences in equation (3). Assume a DM wants to maximize expected value, conditional on her inference on the probability of winning. In the simple choice problems we use this entails a choice rule in which the DM trades off the inference on the log-odds in (3) against the log cost-benefit ratio,  $\ln\left(\frac{c-y}{x-c}\right)$ . For expositional simplicity, we will assume in this section that the log cost-benefit ratio (unlike the log-odds) is objectively perceived, though clearly it will in fact be learned by sampling just as the log-odds are. This assumption will have no impact on our qualitative predictions here but greatly simplifies the exposition.<sup>14</sup> In Appendix A.2, we show that this yields the following *discriminability* equation:

$$\psi_n = \frac{\gamma_n \times \ln\left(\frac{\alpha_n}{\beta_n}\right) - \ln\left(\frac{c-y}{x-c}\right) - \ln(\theta_n)}{\nu_n \times \gamma_n}, \quad (4)$$

where  $\theta_n \triangleq \left(\frac{1-p_0}{p_0}\right)^{1-\gamma_n}$  is the inverse (weighted) prior expectation. This can be interpreted as a measure of “risk aversion” within the model generated by the distorting influence

---

of which apply in the context of our experiments). This happens because the bias introduced in each single inference must be traded off against the resulting reduction in the variance across trials. The estimator used here is optimal in the precise sense that it minimized the mean squared error. Bishop (2006), ch. 3, provides a proof of this optimality in a machine learning context.

<sup>13</sup>Some papers have described recency bias and the importance of the last samples, but as far as we are aware none has truly endogenized the sampling process. E.g., Hau et al. (2008) discuss opportunity costs of sampling and their dependence on the stakes of the experiment. Hertwig & Pleskac (2010) point at the strongly decreasing marginal informational content of additional samples as a possible reason for small samples, without however formalizing this intuition.

<sup>14</sup>It is straightforward to extend the model to include noisy representations of cost-benefit perceptions — see Vieider (2024b). Such noisy representations can, in fact, quantitatively enhance the patterns we describe here in sequence. In our structural model estimates in Online Appendix E, we take explicit account of the effects of sampling on the DM’s beliefs about the cost-benefit ratio.

of the prior.<sup>15</sup> In other words, risk aversion in the model can result from a pessimistic prior expectation by the DM about the types of lotteries she will face.<sup>16</sup> The subscript  $n$  indicates the current sample count, which plays an important role in the characterization of the dynamics underlying the equation.

**Prior expectations and optimal stopping.** Equation (4) trades off two dimensions: the weighted log-odds,  $\gamma_n \times \ln\left(\frac{\alpha_n}{\beta_n}\right)$ , which present evidence in favor of taking the lottery; and the evidence against the lottery, given by the level of pessimism in  $\theta_n$ , and the log-cost benefits. These two dimensions are weighed by their common standard deviation  $\nu_n \gamma_n$ , which measures the degree of confidence in the quantities being traded off. Given the normality assumption on likelihood and prior,  $\psi_n$  will follow a standard normal distribution. The cumulative distribution function of  $\psi_n$  can thus directly be used to predict choice probabilities. In a sampling framework such as DfE, however,  $\alpha_n$ ,  $\beta_n$ , and the derived quantities  $\nu_n = \sqrt{F'(\alpha_n) + F'(\beta_n)}$  and  $\gamma_n = \frac{\sigma^2}{\sigma^2 + \nu_n^2}$ , as well as  $\psi_n$  itself, will all evolve as a function of samples  $n$ . A central intuition underlying our model is that a DM will stop sampling only once she feels that she has sufficient information to reach a decision.

Intuitively, equation (3) thus measures the accumulation of information as samples are taken. If  $\psi_n$  becomes sufficiently positive, indicating information favorable to the lottery, the DM will stop sampling and choose the lottery; if  $\psi_n$  becomes sufficiently negative, the DM will stop sampling and choose the sure amount instead. Unless a threshold is reached, she will keep sampling.<sup>17</sup> Reaching a positive versus negative decision threshold, however, will depend on prior expectations incorporated in  $\theta_n$  (as well as on the log-cost benefit ratio), making the problem asymmetric. Let us assume for simplicity that costs and benefits are equal. Intuitively, a risk averse DM — i.e. a DM with pessimistic prior expectations  $p_0 < 0.5$  — will have less trouble accepting negative evidence (draws of the lower outcome  $y$ ) and reaching the negative discriminability threshold than reaching a positive discriminability threshold after observing the same proportion of draws of the prize  $x$ .

Figure 3 shows an illustration how discriminability  $\psi_n$  evolves. The illustration is based

---

<sup>15</sup>We conceive of the quantities governing the prior  $p_0$ , and  $\sigma^2$ , as constant for the duration of the experiment. This is plausible in our setting where 1) subjects face the same choices in DfE that they have phased in part 1 on DfD; and 2) the experiment is very short.

<sup>16</sup>Note that we do not *assume* the prior to entail risk aversion. We rather treat it as a free parameter through which any underlying risk aversion of the DM may manifest in the model.

<sup>17</sup>What amount of information exactly is deemed ‘sufficient’ by a DM can thereby be subjective and vary from DM to DM. In other words, the precise thresholds used do not affect the qualitative insights derived here. What is important is that the DM will stop sampling once  $\psi_n$  reaches a sufficiently extreme value, passing a subjective discriminability threshold. Note that thresholds may themselves change over time, as is the case in drift-diffusion modeling. Again, this does not affect the qualitative insights we derive here.

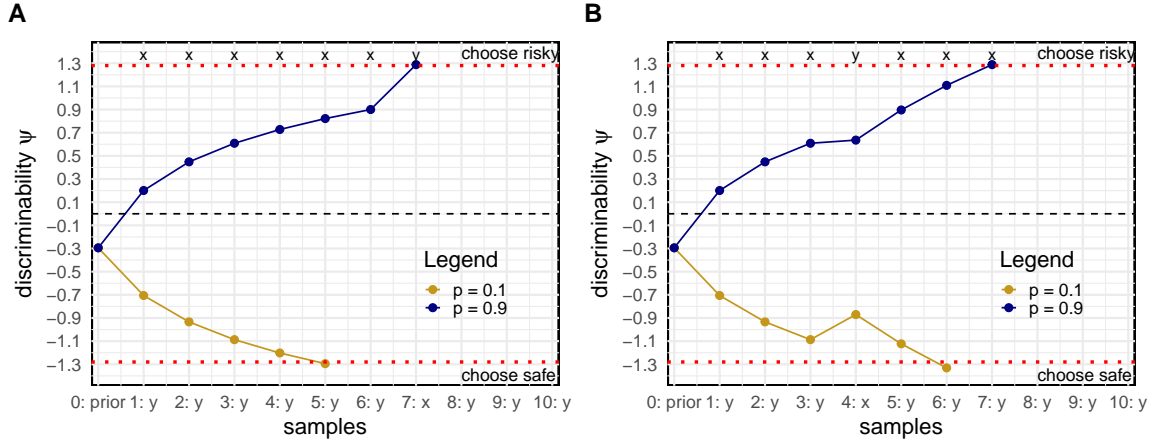


Figure 3: Evolution of discriminability with samples

The figure shows how discriminability in equation (4) evolves with a series of 10 balanced samples. The DM is assumed to be mildly risk averse, with  $p_0 = 0.45$  and to face equal costs and benefits. Panel A illustrates what happens when the outlying observation is drawn in sample 7, whereas panel B illustrates what happens if it is drawn in sample 4 instead.

on a slightly pessimistic DM with  $p_0 = 0.45$ . For simplicity, we assume costs to be equal to benefits, so that their logged ratio drops out of the equation. To focus ideas, we will further assume that in a series of 10 samples, the DM observes exactly one instance of  $x$  if  $p = 0.1$  (and 9 of  $y$ ), and just 1 of  $y$  for  $p = 0.9$ . The discriminability thresholds are set at  $\pm 1.28$ , corresponding to a one-sided test with a 90% confidence level. Panel A depicts the situation in which the less likely outcome is observed in the 7th sample. When sampling from  $p = 0.1$ ,  $\psi_n$  hits the discriminability threshold after a uniform series of 5 ‘failures’  $y$ . The DM stops sampling, and chooses the safe option. While the samples for  $p = 0.9$  are the mirror image of those for  $p = 0.1$ , the evolution of  $\psi_n$  is not. This asymmetry arises from the pessimism in the initial prior. The DM takes 2 more samples, at which point  $\psi_n$  reaches the threshold, and the DM chooses the lottery. Remarkably, this happens even though the 7th and last sample is a draw of a ‘failure’  $y$ .<sup>18</sup> This illustrates the effect of the increase in precision: even though the DM now slightly *under*-estimates the true probability, she nevertheless chooses the lottery because she has high confidence in her estimate.

Panel B illustrates a further implication: that even the position of the outlying observation in the 10 samples will influence the decision when to stop sampling. Here, the less likely event is observed in the 4th draw instead of the 7th. While the positive threshold is

<sup>18</sup>The discriminability equation seems to make a small ‘jump’ upon sampling of the failure  $y$ . This arises from the definition of sampling variance  $\nu_n^2 = F'(\alpha_n) + F'(\beta_n)$ , since the marginal increase of the trigamma function is largest at small values.

still reached after 7 samples, reaching the negative one now requires one more sample than before.<sup>19</sup> There is also a general implication that emerges from this illustration: all sequences of 10 samples we used in the example are accurate in the sense that taking 10 samples results in the true underlying choice proportion being sampled. Nevertheless, the endogenous decision on when to stop sampling produces sampling error in all four cases. This illustrates how sampling variance and sampling error interact, given that the former — in combination with the prior expectation — will determine when a DM stops sampling. In reality, 10 samples will typically *not* correctly reflect the true probability, as illustrated in figure 2 above. This will further increase the distortions introduced by the endogenous sampling-stopping decision.

The interaction between precision and sampling error results in a testable insight: risk averse DMs should sample more from large probability lotteries than from small probability lotteries on average, whereas risk-loving or optimistic DMs should do exactly the opposite. This prediction is new, and has not previously been examined in the literature, making it a test that is diagnostic of the value of our model for the prediction of sampling behavior.

**Implications for risk-taking.** The decision on when to stop sampling, described above, will in turn contribute to determining risk-taking patterns (jointly with the log-cost benefits). Risk averse DMs facing small probability lotteries will stop sampling early because the accumulation of failures arising from sampling error reinforces the prior expectation. This will suggest a choice of the sure option, explaining widespread risk aversion for small probabilities. For large probabilities, however, the initial series of successes drawn by a majority of DMs clashes with their pessimistic expectations. This prevents the positive discriminability threshold from being reached, and leads DMs to take larger samples. These larger samples will still on average look favorable due to sampling error, which we have seen to only decrease very slowly in the number of samples (cfr. figure 2, panel A).<sup>20</sup> The increased precision of the larger samples will concomitantly reduce the weight put on the pessimistic prior expectation (i.e.,  $\theta = \left(\frac{1-p_0}{p_0}\right)^{1-\gamma}$  converges to 1 as  $\gamma$  increases, and  $\ln(\theta)$  goes to 0). This explains relatively high risk-taking for large probabilities, and thus positive probability-dependence of risk-taking in DfE. It also results in a testable prediction: the

<sup>19</sup>An interesting consequence of this sort of stopping decisions may be apparent recency effects in DfE, as discussed e.g. by Erev & Barron (2005). In the context of our model, however, such recency effects would be mostly driven by the decision on when to stop sampling (see also Wulff et al. 2018, for a discussion of this point).

<sup>20</sup>It is important to note that these are *average* patterns. Some DMs may draw very favorable samples from small probability lotteries or very unfavorable samples from large probability lotteries, and thus take the opposite decisions.

substantial risk-taking for large probabilities — and the positive probability-dependence more generally — ought to be driven primarily by DMs who are ex ante highly risk averse.

### 3 Sampling in Decisions from Experience

#### 3.1 Experiment description

In Experiment 1, we replicate the GAP, as shown in figure 1 in the introduction. Subjects face 18 distinct binary choices between a sure amount  $c$  and a lottery paying  $x > c$  with probability  $p$ , or else  $y = 0$ . We further randomly pick 4 choice problems to be repeated. The lotteries vary  $p$  across 0.1, 0.15, 0.2, 0.8, 0.85 and 0.9 and vary payoffs  $x$  and  $c$ . The sure amounts  $c$  for a given probability include the expected value ( $EV$ ) of the lottery, and two amounts that are symmetric around the EV of the lottery (i.e.  $c = EV(x, p) \pm h$ , where  $h$  is \$0.3 or \$0.4). This will allow us to get a rich picture of behavior, and is crucial to identify probability-dependence in risk-taking. We did not include intermediate probabilities because we followed the DfE literature in our task selection, where the use of intermediate probability is rare since they are typically not very informative for the GAP.<sup>21</sup> Lotteries are described to subjects as “bags,” containing 20 “coins,” each of which is worth a different amount of money. At the end of the experiment, a lottery is randomly selected and a single coin is drawn from the bag to determine the subject’s payment.

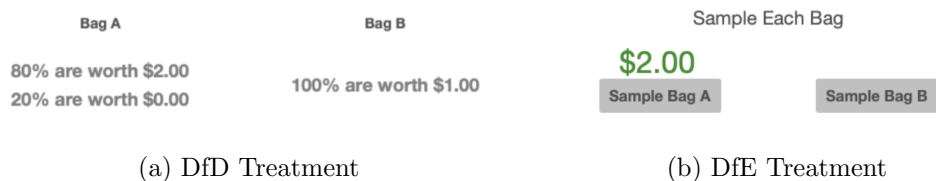


Figure 4: Screenshots from Experiment 1.

**Treatments.** Experiment 1 consists of two treatments. In the DfD treatment, the subject is explicitly told the properties of each lottery (i.e., the contents of each bag); Figure 4a shows a screenshot. A pair of radio buttons below the lottery description allows the subject to make and submit a choice between the two lotteries. In the DfE treatment, the subject is instead shown two buttons, one for each of the two lotteries/bags. Figure 4b shows a screenshot. When the subject clicks on the button, she is shown a single realization of the lottery (i.e., a single draw from the bag, *with* replacement). The subject is told nothing

<sup>21</sup>This was also meant to limit the number of clicking necessary in the experiment, given that the forced sampling treatments described below require 41 mouse clicks per task.



about either lottery and must learn all of their properties by sampling. The subject in the Figure 4b example has just clicked the “Sample Bag A” button and drawn \$2. Each sample is shown for 0.5 seconds. Subjects are allowed to sample as many (or as few) times as they like from the two bags, with no time constraints. Below the sampling buttons are the same two radio buttons shown beneath the lottery descriptions in DfD, and the subject can choose one of the lotteries to determine her payment whenever she is ready.

**Stages.** Each session in the experiment proceeds in two stages. In Stage 2, subjects experience their main treatment: 22 randomly ordered lottery choices under DfD or DfE, depending on treatment. In Stage 1, subjects face the same 22 binary choice tasks under DfD (in a different random order). We included Stage 1 for several reasons. First, doing this allows us to examine the GAP both within-subject (by comparing Stage 1 and Stage 2 in the DfE treatment) and between-subjects (by comparing Stage 2 in the DfE vs. DfD treatment). Second, including Stage 1 is useful for fixing prior beliefs about lotteries and linking DfD and DfE behavior.

**Implementation.** We ran 99 subjects through the DfD treatment and 99 subjects through the DfE treatment on Prolific. We paid all subjects \$6 and selected 10% of them to be paid based on a lottery outcome from a randomly selected task. The median subject spent 18 minutes in the experiment and the average subject earned \$18.67 per hour. Instructions, including 4 comprehension questions, are included in Online Appendix F.

### 3.2 Results: Sampling, Sampling Error, and Risk-Taking

**Sampling patterns.** We start by testing the key edogenous sampling predictions: (i) sampling behavior should vary with the probability of the prize,  $p$ ; and (ii) this dependence should vary according to the subject’s pre-existing level of risk aversion as captured by the prior mean. To test this, in first approximation we categorize DfE subjects according to their risk aversion using their propensity to choose risky lotteries in Stage 1 by quantifying the proportion of Stage 1 risk averse choices subjects make.<sup>22</sup>

In Figure 5, panel A, we plot the mean number of samples taken from the risky option in Stage 2 as a function of probability  $p$  for subject classified as High and Low risk aversion based on a median split of risk averse choices in the first, DfD stage of the experiment. We find clear evidence of the predicted pattern: highly risk averse subjects sample substantially

---

<sup>22</sup>As we will explain shortly, we expect behavior in DfD to also be affected by sampling variance, so that this measure is only a *proxy* for risk aversion as captured by the prior. The results we report are, however, robust to using the structurally estimated prior mean instead.

more at high than at low probabilities; relatively risk tolerant (Low risk aversion) subjects show (somewhat weaker) evidence of the reverse sampling pattern. Given that most subjects in our sample are risk averse this results in an overall average tendency for subjects to sample more for larger than smaller probabilities. These results strongly support the idea that sampling precision interacts with prior beliefs to determine sampling behavior in DfE.

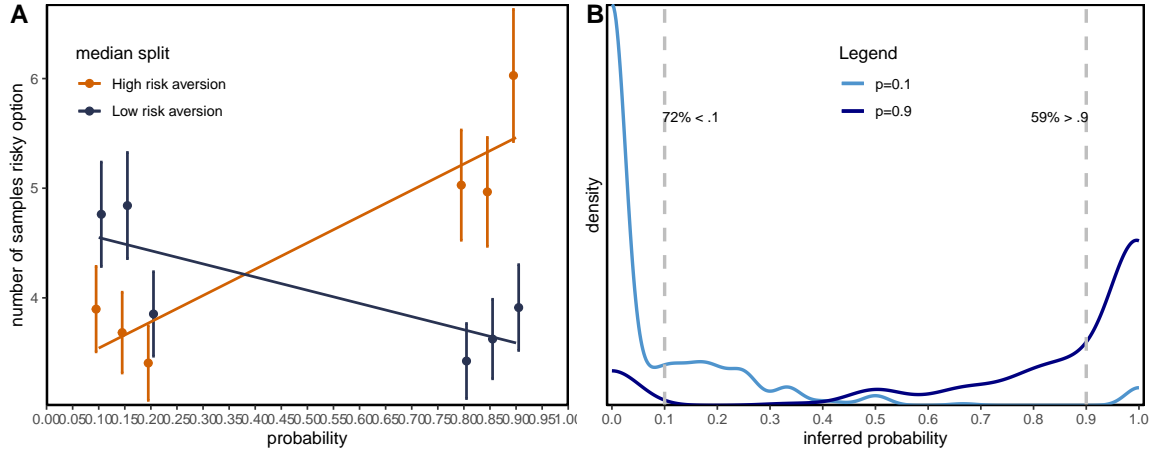


Figure 5: Samples by probability and risk aversion

Panel A shows the number of samples taken from the risky option by probability and risk aversion. Risk aversion is assessed as the proportion of safe choices in the first, DfD part of the experiment, after removing repeated tasks. The categorization is obtained using a median split. Error bars show  $\pm 1$  standard error. Panel B shows the distribution of actual sampling error in the samples taken for a small probability of  $p = 0.1$  versus a large probability of  $p = 0.9$ . Dashed vertical lines show the true underlying probabilities generating the samples.

**Endogenous sampling choices drive sampling error.** Panel B in figure 5 shows the resulting samples drawn for a small probability of  $p = 0.1$  and for a large probability of  $p = 0.9$  (findings for other probabilities are similar, and shown in Online Appendix E). The figure shows the direct consequence of smaller samples taken for  $p = 0.1$  on average by a risk averse subject population: the error of underestimating a probability of  $p = 0.1$  is clearly more frequent than the error of over-estimating  $p = 0.9$ . Across all small probability lotteries, our subjects gather samples that produce a smaller probability than the true one in 66% of cases overall, and an accurate sample in 3.4% of cases. For large probability lotteries this is reversed, with 55% of samples over-estimating the true probability, and only 2.2% resulting in a correct estimate. Sampling error is clearly more severe for small probabilities than for large probabilities — a direct consequence of the smaller number of samples taken for small probabilities on average.

**Individual-level analysis.** So far we have only shown *aggregate* patterns. An important conclusion from our simulations was, however, that behavior will depend on individual-level

expectations, as well as on the fortuitous samples drawn — including both the proportion of winning outcomes, and the sequence in which these proportions are drawn. Here, we use regression analysis to show that the results we document above are robust (i) to using continuous measures of risk aversion; (ii) to conducting this analysis at the individual level; and (iii) to using structurally estimated measures of prior expectations instead of the proxy constituted by proportion of choices of the sure amount in DfD. All of these effects are fully in line with the model predictions delineated above, and thus support our account of endogenous sampling, and how it leads to sampling error.

| dep. var:              | number of samples |                  |                  | abs. sampling error |                   |                   |
|------------------------|-------------------|------------------|------------------|---------------------|-------------------|-------------------|
|                        | reg. (1)          | reg. (2)         | reg. (3)         | reg. (4)            | reg. (5)          | reg. (6)          |
| probability            | 0.288<br>(0.082)  | 0.331<br>(0.085) | 0.326<br>(0.083) | -0.027<br>(0.008)   | -0.026<br>(0.008) | -0.020<br>(0.008) |
| risk aversion          |                   | 0.569<br>(0.291) | 0.562<br>(0.266) | -0.016<br>(0.012)   | -0.021<br>(0.011) | -0.009<br>(0.010) |
| prob $\times$ risk av. |                   | 0.618<br>(0.087) | 0.645<br>(0.086) | -0.025<br>(0.009)   | -0.030<br>(0.008) | -0.004<br>(0.001) |
| samples                |                   |                  |                  |                     |                   | -0.018<br>(0.009) |
| constant               | 3.635<br>(0.283)  | 3.705<br>(0.289) | 3.700<br>(0.284) | 0.046<br>(0.011)    | 0.006<br>(0.023)  | 0.072<br>(0.011)  |
| observations           | 2178              | 2178             | 2178             | 2178                | 2178              | 2178              |
| subjects (clusters)    | 99                | 99               | 99               | 99                  | 99                | 99                |

Table 1: Regression analysis of samples taken and sampling error

Regressions in the table are based on a Bayesian outlier-robust regression model. Robust regression is implemented by means of a student-t distribution with 2 degrees of freedom, with random intercepts to cluster errors at the subject level. Regressions (1), (2), and (3) use the total number of samples from the risky option as dependent variable. Regressions (4), (5), and (6) use the absolute sampling error, defined as the true probability minus the inferred probability for small probability lotteries, and as the inferred probability minus the true probability for large probability lotteries, as dependent variable. Numbers in parentheses indicate standard errors. Risk aversion is captured by the proportion of risk averse choice in phase 1 DfD in columns (1), (2), and (4), and by the inverse log-odds prior  $\ln\left(\frac{1-p_0}{p_0}\right)$  in regressions (3), (5) and (6). Probability and risk aversion are normalized by taking z-scores.

Table 1 shows regressions detailing individual-level patterns. Regression (1) shows that samples taken increase in the probability of winning across all subjects. Regression (2) uses proportions of risk averse choices in the initial DfD phase to show that the larger overall samples are mainly driven by risk aversion, and that probability-dependence of the number of samples taken strongly increases in pre-existing risk aversion. Regression (3) further probes the robustness of these results by instead using the theoretically correct measure of

pessimism in the prior expectation,  $\ln\left(\frac{1-p_0}{p_0}\right)$ , which we obtain from structural estimations of equation (4) from the first phase DfD data.<sup>23</sup> All the results remain stable.

Regressions (4) through (6) present regressions of the absolute sampling error (coded in the direction of underestimating the small likelihood event). Regression IV uses the independent variables from regression (2) to show that (i) sampling error decreases in the probability of winning for average levels of risk aversion; and (ii) that sampling error for large probabilities decreases (and sampling error for small probabilities *increases*) in the level of pre-existing risk aversion. Regression (5) shows that all of these effects are stable to using the correct prior expectation from regression (3) instead. Finally, regression (6) adds the number of samples as an explanatory variable, and shows that sampling error indeed decreases in the samples taken, as one would expect. The number of samples taken thereby absorbs almost the entire effect of pre-existing risk aversion in the reduced form equations (4) and (5). Taken together, these results show that 1) the number of samples taken increases in the level of pre-existing risk aversion; 2) risk aversion particularly increases samples for large probabilities, while it lowers them for small probabilities; and 3) the same individual characteristics also determine the extent of sampling error. This constitutes a first key piece of evidence in support of the mechanisms predicted by our noisy cognition model.

**Number of samples, sampling error, and risk-taking.** To complete the picture of what drives behavior in DfE, we next look at risk-taking choices. Table 2 reports a series of Probit regressions to investigate drivers of risk-taking at the individual level. Regressions (1) and (2) show the reduced form regressions using the same characteristics used above to predict sampling behavior and absolute sampling error (using the risk averse choice proportion and the estimated inverse prior log-odds, respectively). Risk taking increases in the probability of winning. Remarkably, however, risk taking particularly increases in the probability of winning for DMs *with the highest pre-existing degree of risk aversion*. This result constitutes direct evidence for the mechanism predicted by our model.

Regression (3) further adds the number of samples from the risky option and the sampling error (coded as the error in samples in favor of the lottery). Both are highly significant predictors of the level of risk-taking. They also takes up most of the effect previously captured by the interaction between the probability of winning and pre-existing risk aversion, which

---

<sup>23</sup>Note that to be applicable to the DfD data, the true underlying log-odds  $\ln\left(\frac{p}{1-p}\right)$  have to be substituted for  $\ln\left(\frac{\alpha}{\beta}\right)$  in that equation. Section A in the Online Appendix discusses the theoretical rationale for this substitution, and section E provides the details about the econometric estimation.

| dep. var:                    | choice of lottery over sure amount |                  |                   |                   |
|------------------------------|------------------------------------|------------------|-------------------|-------------------|
|                              | reg. (1)                           | reg. (2)         | reg. (3)          | reg. (4)          |
| probability                  | 1.134<br>(0.095)                   | 1.199<br>(0.095) | 1.775<br>(0.123)  | 1.714<br>(0.131)  |
| risk aversion                | 0.053<br>(0.105)                   | 0.112<br>(0.110) | -0.188<br>(0.151) | -0.162<br>(0.147) |
| prob $\times$ risk av.       | 0.625<br>(0.093)                   | 0.735<br>(0.099) | 0.259<br>(0.050)  | 0.029<br>(0.056)  |
| nr. of samples               |                                    |                  | 0.877<br>(0.052)  | 1.534<br>(0.085)  |
| samp. error for lottery      |                                    |                  | 0.673<br>(0.116)  | 0.790<br>(0.119)  |
| samp. error $\times$ samples |                                    |                  |                   | 1.160<br>(0.017)  |
| constant                     | 0.398<br>(0.106)                   | 0.436<br>(0.104) | 0.706<br>(0.147)  | 0.373<br>(0.153)  |
| observations                 | 2178                               | 2178             | 2178              | 2178              |
| subjects (clusters)          | 99                                 | 99               | 99                | 99                |

Table 2: Regression analysis of risk-taking

The table shows Bayesian Probit regressions of risk-taking on a number of independent variables. Errors are clustered at the subject level using random intercepts. Probability, risk aversion, sampling error, and samples from the risky option are normalized by taking z-scores. The sampling error is defined in the direction of favoring the lottery. Standard errors are shown in parentheses.

nevertheless remains significant. Regression (4) further adds the interaction between sampling error the number of samples taken. Risk-taking strongly increases in this interaction. At the same time, the interaction between a pessimistic prior expectation and the probability loses its significance. This highlights the interactive role sampling error and number of samples play in our model: sampling error in favor of the lottery will be most influential in determining decisions when the DM has high confidence in the sampled proportion.

Taken together, the regressions above strongly support the mechanism predicted by our model to drive positive probability-dependence of risk-taking in DfE. Regressions (1) through (3) in table 1 illustrate the effect of pre-existing risk aversion, and its interactions with probability, on the number of samples taken. Regressions (4) through (6) in that same table illustrate the effect this has on sampling error. Finally, regressions (1) through (4) in table 2 illustrate the effect that pre-existing risk aversion interacting with the probability of winning — mediated by the number of samples, the sampling error, and their interaction — have on risky choice.

### 3.3 Experiment 2: DfE+forced Treatment

In Experiment 2, we attempt to eliminate the GAP by forcing DfE subjects to sample from each lottery (i) using a representative sample and (ii) via a relatively large number of draws. The balanced nature of the samples is meant to eliminate sampling error. The large number of draws is meant to increase sampling precision, thereby inducing subjects to rely on the sampled information more than they rely on their prior expectation.



Figure 6: Screenshot from the DfE+forced treatment (Experiment 2).

We do this using the DfE+forced treatment, pictured in Figure 6. This treatment is identical to DfE except that subjects are required to sample all twenty “coins” from each bag (lottery) *without replacement* before making a choice between lotteries. Below each button, the subject is shown how many times she has sampled from each bag and the total number of draws she must make in total (set to 20 in this treatment). The radio buttons for submitting the final lottery choice do not appear on the subject’s screen until she has sampled all 20 coins from each bag. In terms of the model, requiring the subject to exhaustively sample a frequentist representation of each lottery means that subjects observe samples  $\alpha, \beta$  for each lottery such that  $\frac{\alpha}{\alpha+\beta} = p$ , removing scope for sampling error. By setting the number of elements in the frequentist representation to 20, we force subjects to sample far more times than they are observed to do in the DfE treatment, thereby increasing the precision  $\nu^{-2}$ . In all other respects the experiment is identical to the DfD and DfE treatment.

**Forced sampling eliminates probability-dependence in DfE.** Figure 7 plots choice behavior from DfE+forced, and reproduces behavior from DfE for comparison. As predicted, forced sampling produces a dramatic effect on behavior, particularly in reducing the high levels of risk taking observed for large probabilities. Importantly, as predicted, DfE+forced does this largely by eliminating probability-dependence.

To analyze this more systematically, we calculate choice proportions and their standard errors for each of the 18 tasks. We then aggregate choice proportions across tasks weighing them by the inverse of their squared standard errors, as done in meta-analysis or measurement error models. Regressing the choice proportions of the lottery on the probability of winning provides a direct test of probability-dependence in the choice proportions (see

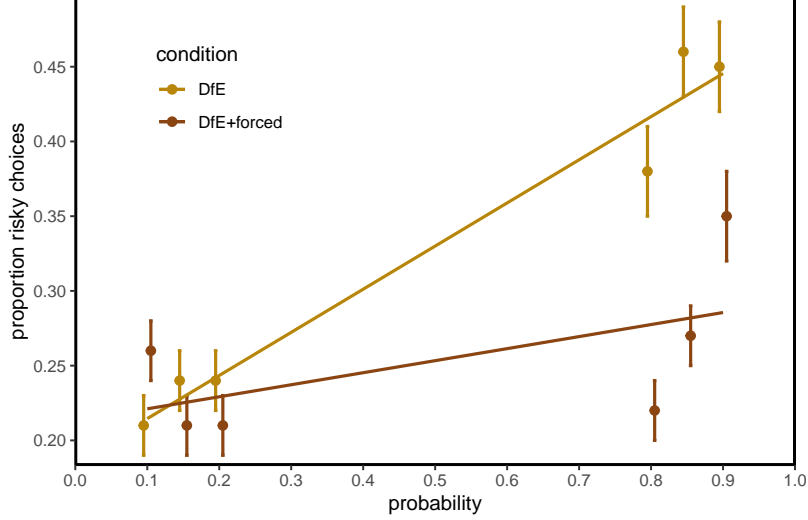


Figure 7: Effects of forced sampling in DfE on choice proportions

The figure shows the effects of forced complete sampling in DfE from both options. The error bars indicate  $\pm 1$  standard error.

Online Appendix C for details). In DfE, this produces a coefficient on  $p$  of 0.284, with a credible interval of  $[0.226, 0.347]$ , showing how risk-taking systematically increases in the probability of winning. Regressing choice proportions observed after forced sampling on the probability of winning the prize, we find a slope of 0.080, with a CrI of  $[-0.020, 0.182]$ . This slope is significantly smaller than in DfE. The slope is also not significantly different from 0 — positive probability-dependence of risk-taking in DfE has disappeared upon forced sampling.<sup>24</sup> Notwithstanding these significant changes on the DfE side, however, the GAP does not close: although it has narrowed to 8.4 pp, it remains substantial as well as statistically significant, with a CrI of  $[0.053, 0.115]$ .

Our results are consistent with previous findings on forced sampling in DfE. Pioneering the use of forced sampling in DfE, Ungemach et al. (2009) found the GAP to narrow, but not close (see also Cubitt et al. 2022). Our conclusions are fully consistent with this finding, but strengthen it further. In particular, our richer test stimuli allowed us to test probability-dependence after forced sampling directly, and to show that it disappears — something Ungemach et al. (2009) could not do due to the smaller number of task and absence of variation in the EVs of the choice options conditional on a given probability.

<sup>24</sup>One reason why there is still a slight positive tendency in risk taking could be memory effects. While the previous literature using similar settings to our own has not found much of a role for memory (Ungemach et al. 2009, Cubitt et al. 2022), we do find that the more recent half of samples has a slightly stronger effect on risk-taking than the first half. This is consistent with effects of memory documented by Bohren et al. (2024), albeit in a somewhat different setting.

## 4 Closing the Gap: Forced sampling in DfD

Forced sampling leads to dramatic changes to DfE behavior, but it does not close the GAP. This happens because — while forced sampling removes probability-dependence in DfE — DfD still exhibits risk-taking that declines in the probability of winning. The key reason for this emerging from our model is that very similar factors also afflict DfD. Risk-taking decreasing in the probability of winning in DfD, we hypothesize, is a consequence of imprecision in the perception of probabilities in DfD. Because of this, in order to fully close the GAP, we have to increase the precision in probability perceptions in DfD in a manner symmetric to the way we did in DfE.

We start from the observation that probabilities will have to be neurally represented by spikes or action potentials before entering the decision processes. This will result in a noisy signal  $r$  for the described log-odds, as modelled by Khaw et al. (2025) and Vieider (2024b). Here, we hypothesize that this noisy signal can be conceived of as a ratio of “neurally sampled” evidence in favor of the lottery and against it, and summarized by a quantity  $\ln(\hat{\alpha}/\hat{\beta})$  similar to the ratio of real samples used to characterize DfE above. This follows the seminal discussion of log-odds coding by Gold & Shadlen (2001), who forcefully argue that it is efficient for the brain to summarize evidence about an uncertain hypothesis using a population of neurons in favor of the hypothesis and a population of “anti-neurons” summarizing the evidence against.

This conceptual framework yields several insights: 1) the noisy signal for the true log odds can be conceived of as a ratio of firing rates signaling evidence in favor and against the lottery,  $\ln\left(\frac{\hat{\alpha}}{\hat{\beta}}\right) \sim \mathcal{N}\left(\ln\left(\frac{p}{1-p}\right), \nu^2\right)$ ; 2) finite neural spike counts or activation potentials making up  $\hat{\alpha}$  and  $\hat{\beta}$  will yield noisy representations in single trials; and 3) even though the representation will be correct *on average*, thus avoiding the sampling error that affected DfE, the lack of precision in the signals will yield regression to the mean of the prior just as in equation (3) (see online appendix A for further details).

**Closing the Gap.** This explanation suggests a distinctive test for the hypothesis that imprecision in neural representations is responsible for classic probability-dependence in DfD. If it is true that probability-dependence in DfD is driven by the noisy perception of described probabilities as hypothesized above, then providing additional information in the form of forced, balanced samples ought to remove this probability-dependence, just as it did in DfE. This makes for a powerful test because 1) from the point of view of standard models such as prospect theory, the information gleaned from samples is *fully redundant*



in the DfD treatment; and 2) combining sampling information with described information ought to allow us to increase the precision of the neural log-odds representations, thus reducing or even eliminating probability-dependence in DfD. This provides a particularly crisp test of noisy coding accounts of probability-dependence, since we can directly act on the supposed causes of the phenomenon to try and remove it.

### 4.1 Experiment 3: the DfD+Forced Sampling Treatment

In Experiment 3, we attempt to eliminate the GAP by forcing DfD subjects to redundantly sample large, representative samples from each lottery. In DfD+forced we show subjects the same information about lotteries as we do in the DfD treatment (pictured in Figure 4), but we also provide subjects the sampling tools pictured in Figure 6 below the explicit description, and force subjects to draw 20 times from each just as in DfE+forced. Indeed, the DfD+forced treatment is identical to the DfE+forced treatment, except that lotteries are fully described to the subject prior to, during and after sampling. This is an original treatment that has not been tested before in the literature.

**Forced sampling eliminates probability-dependence from DfD.** Panel A of Figure 8 shows the effect of forced sampling in DfD, by plotting average choice proportions for DfD+forced and (for comparison) DfD. As predicted, we find that forced sampling has *exactly the reverse effect* on DfD as on DfE. At small probabilities, we find a sizeable *decrease* in risk taking at most probabilities.<sup>25</sup> For large probabilities, on the other hand, risk taking *increases* with forced sampling in DfD. Thus, just as predicted by noisy cognition models like ours, providing completely redundant information to subjects has a sizable effect on choices in DfD. While probability-dependence in DfD is  $-0.172$ , CrI  $[-0.239, -0.102]$ , it has disappeared after forced sampling (slope  $0.058$ , CrI  $[-0.039, 0.157]$ ).

**Closing the GAP.** What does the elimination of sampling error and the increase in precision via forced sampling (in both DfE and DfD) do to the GAP? Panel B of figure 8 shows that the choice proportions are now very similar. Probability-dependence in DfD+forced, at  $0.058$  (CrI  $[-0.039, 0.157]$ ), and in DfE+forced, at  $0.080$  (CrI  $[-0.020, 0.182]$ ), are not significantly different from each other. Nor does risk-taking in either treatment show any sign of probability-dependence: the opposite patterns in DfD and in DfE disappear upon forced sampling, converging to mild risk aversion as in standard EUT. As our structural

---

<sup>25</sup>The exception is  $p = 0.15$ . This is, however, in part caused by the aggregation across different values of sure payments,  $c$ . For this particular probability, the changes across different sure amounts go in opposite directions canceling each other out – see Online Appendix E for the plot broken down by values of  $c$ .

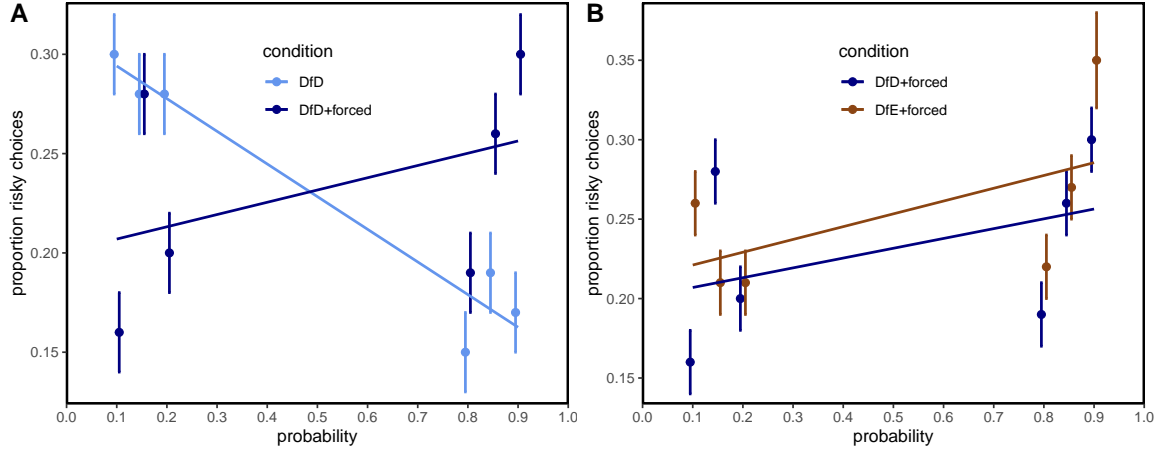


Figure 8: Effects of forced sampling on choice proportions

The figure shows the effects of forced complete sampling in description-based choice. Panel A shows the effect of forced sampling in DfD+forced on choice proportions for different probabilities, and compares it to DfD. Panel B directly juxtaposes choice proportions in DfE+forced and DfD+forced. The lines are fitted by linear regression to the average choice proportions by probability. The error bars indicate  $\pm 1$  standard error.

estimations in Online Appendix E show, subjects furthermore become much more similar to each other, and the absence of probability-dependence becomes the rule in the data.

To provide a more nuanced picture, and to examine the GAP directly, figure 9 shows differences in choice proportions between DfD and DfE for all 18 tasks. Panel A shows the original GAP between DfD and DfE. We use a measure  $g$  capturing the difference in choice proportions, defined so that positive values correspond to behavior typically documented in the literature for the standard GAP — more risk-taking in DfD than DfE for small probabilities, more risk taking in DfE than DfD for large probabilities. We then meta-analytically aggregate the GAP across tasks, which yields an estimate of the overall GAP as well as correcting for random sampling variation in single tasks (details in Online Appendix C). In the *absence* of forced sampling, the GAP is significant in 12 out of 18 tasks when looking at the raw choice proportions, and in 13 out of 18 tasks in the meta-analytic posterior.<sup>26</sup> At 15.7 percentage points ( $pp$ ), with a 95% credible interval of  $[9.7, 21.8]$   $pp$ , the GAP is significant and large measured against the meta-analytic average reported by Wulff et al. (2018), which comes to book at 9.7  $pp$ .

Panel B compares description-based and experience-based choice proportions after forced sampling (DfE+forced vs. DfD+forced), and shows that the GAP disappears in these treatments. We find no significant gap for *any* of the 18 choice proportions in the meta-

<sup>26</sup>The exceptions in which the GAP is not statistically significant at conventional levels are small probability tasks with  $c \geq px$ .

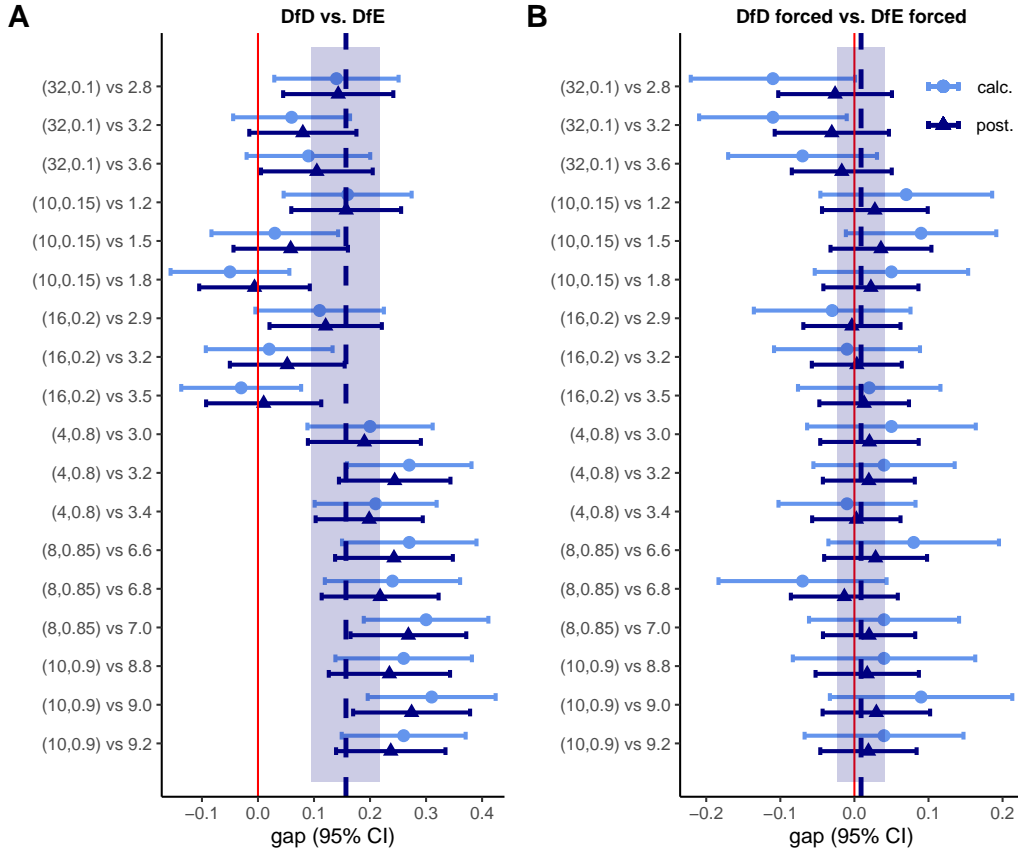


Figure 9: Meta-analysis of the GAP

Panel A shows a forest plot of the gap for our standard implementation of DfD versus DfE. Panel B shows a forest plot of the GAP after forced sampling both from description and from experience. The light blue circles, labeled ‘calc.’, indicate the raw differences in choice proportions in the data,  $g$ . The dark blue triangles, labeled ‘post.’, indicate the inferred posterior parameters,  $\hat{g}$ . The thick, dashed vertical line indicates the meta-analytic posterior mean,  $\omega$ , and the shaded rectangle indicates the 95% credible interval around that estimate.

analytic posterior. In the one case in which we see a significant gap in the raw choice proportions, this gap goes in the *opposite* direction of the standard GAP. At 0.9 pp (95% credible interval of  $[-2.3, 4.1]$  pp), the meta-analytic posterior mean is arbitrarily close to 0. The GAP has closed. Our results thus provide strong evidence that the decision-experience GAP is a consequence of the two elements we expect forced sampling to remove.

## 5 Free Sampling from Described Choice Options

Representations of lotteries in DfD are noisy, indicating some natural limit to the precision with which probabilities can be mentally represented. This raises two intriguing questions: First, will subjects sample when given fully described options, even when they are not forced

to do so? Second, if subjects do indeed sample, will the sampling error flip the probability-dependence in risk-taking to qualitatively resemble the pattern observed in DfE? Assuming that subjects combine the sampled information — which will inevitably be affected by error — with the unbiased description, we may indeed expect probability-dependence of risk-taking to change from negative in DfD to positive in DfD+free.

In the DfD+free treatment, we show subjects the same information about lotteries as we do in the DfD treatment, but we also provide subjects the sampling tools just like in the DfD+forced treatment. Other than in DfD+forced, however, the radio buttons to indicate a choice appear from the very start. Subjects are told explicitly that they can sample if they want to but that they do not have to, and that they can also indicate their decision directly without sampling. We ran this treatment on Prolific with 101 subjects using otherwise identical tasks and procedures as in DfD.

We do indeed observe that subjects sample when given the chance to do so, suggesting at least partial awareness of coding noise in DfD. Across subjects and tasks, the average number of samples is 1.74, of which 1.4 are taken from the risky option. Only 8 out of 100 subjects never sample at all, but most subjects take relatively few samples. Samples are highest at the beginning, with 4.7 samples being taken on average across all subjects in the first round. This declines rapidly to some 2.9 samples on average in the second round, and to 2.4 in the third. After round 8, the average settles to a steady level of 1.3 samples per task and subject. The fact that DMs do sample fully redundant information seems remarkable in our context, given the high opportunity costs of subjects on Prolific.

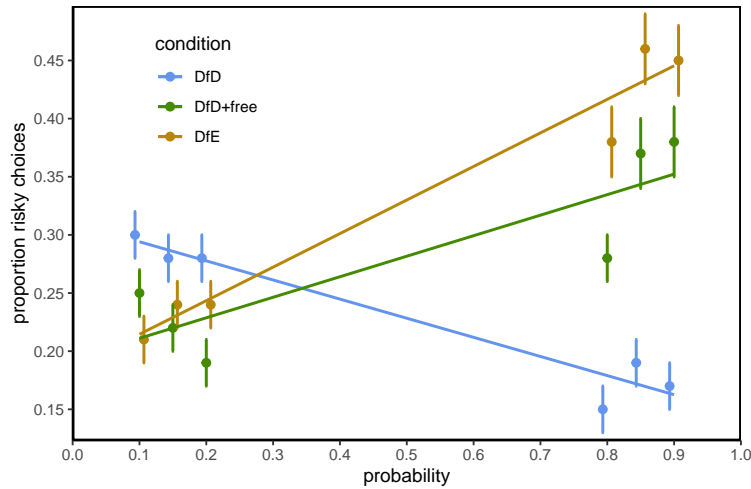


Figure 10: Structural estimates, DfD versus DfE

We next examine what happens to choice behavior once free sampling is introduced. Our model raises an intriguing question: may free sampling in DfD introduce sampling error into DfD? The question arises simply because, although subjects are given an objective description of the probabilities, our model suggests that actual samples drawn are combined with the unbiased neural samples representing the evidence in favor and against the lottery. Small samples, however, will suffer from the same issues we have seen in DfE: they will tend to under-estimate the likelihood of observing the rare event, so that our model predicts that they will yield biased updates of the true log-odds.

Figure 10 shows the raw choice proportions in DfD+free, and directly juxtaposes them with the choice proportions in DfD and in DfE. The difference from DfD is very large, with somewhat less risk-taking for small probabilities, and much more risk-taking for large probabilities. This results in a positive dependence of choice proportions on the probability of winning, with a slope of 0.174, and a 95% credible interval of [0.141, 0.241]. This suggests that sampling error indeed affects DfD+free choices, just like predicted by our model. The effect is indeed strong enough to considerably narrow the GAP in the opposite direction when examining it meta-analytically: at 3.5 pp, the average GAP is now small, and (just) not statistically significant, with a 95% CrI of  $[-0.001, 0.072]$ .

The picture is more nuanced when looking at probability-dependence directly. Although we now see positive probability-dependence of risk-taking in DfD+free, the influence of the description is strong enough to keep the probability-dependence significantly smaller than observed under DfE (where we have observed a slope 0.284, with a credible interval of [0.226, 0.347]). This shows that samples from description — while affecting probabilities in the way predicted by our model — are still balanced against the description provided on the screen, with choices indicating an aggregation of the two types of information.

Our findings are consistent with studies that have investigated the effect of providing feedback after payoff-relevant choices under the form of single draws from the chosen option. Van de Kuilen (2009) studied the effect of feedback provision after risky choices in a prospect theory framework. He concluded that providing feedback shifted behavior towards linearity in probability weighting, but could not fully test this proposition due to exclusive focus on probabilities  $\geq 0.5$  (see also van de Kuilen & Wakker 2006). Jessup et al. (2008) and Tymula et al. (2023), who use both monkeys and humans as subjects, provided feedback after choices for a large number of trials. While none of these studies focuses on the GAP, they show that classic probability-dependence in DfD reverses upon the provision of feedback,

resulting in the type of positive probability-dependence observed in DfE.

Our model sheds new light on these findings: unless independently and identically distributed samples are extremely large, they will introduce sampling error into DfD. The reason this happens lies in the imprecision in the mental representations of probability that are explicitly described: the residual uncertainty in their mental representations provides opportunities for additional information to affect those representations. The feat of closing the GAP with DfE by acting on DfD is remarkable inasmuch it achieves something that acting on DfE alone has never achieved — it closes the GAP by manipulating one of the two sides only. Guided by our model, we introduced sampling error into DfD, all the while keeping precision relatively low due to the few samples added. This dramatically narrows the GAP when people can sample freely, with DfD+free approaching the type of positive probability-dependence characteristic of DfE.

## 6 Discussion

In this paper we show that probability-dependent risk-taking and the description-experience gap — two key phenomena in the lottery choice literature — are a consequence of the incomplete and imprecise ways decision makers perceive and represent information. Reducing the imprecision of subjects’ beliefs by forcing them to observe redundant information causes probability-dependence in risk-taking to disappear and closes the description-experience gap. In addition to shedding significant light on a key mystery in the literature, we believe there are several broader implications of our findings.

First, our results show the reach of the noisy cognition approach by extending it from description-based choice to experience-based choice. Noisy cognition thereby organizes a key paradox under existing descriptive models of choice — the description experience gap — showing the added value of the approach over existing models. Our sampling-based characterization furthermore allows us to present a particularly crisp test of standard probability-dependence in risk-taking when choice options are described: by forcing subjects to take large, balanced samples from both choice options, we manage to completely eliminate probability-dependence in risky choice. This treatment effect is difficult to account for via alternative explanations that are not similarly rooted in cognitive imprecision. It further shows that probability-dependence cannot be attributed to preferences, but is purely an outgrowth of noisy cognition — a conclusion that goes beyond what has been shown in the previous noisy cognition literature.

Noisy coding models hypothesize that descriptive failures of benchmark models like expected utility theory (von Neumann & Morgenstern 1944, Savage 1954) are a consequence, not of non-standard preferences, but rather of what we have called imprecision in mental representations, driven by limitations in the way the brain encodes information. This raises the question of what type of behavior a welfare-maximizing policy maker should take into account. By strongly reducing imprecisions in mental representations, our treatment interventions reveal mild, probability-independent risk aversion as a candidate for welfare-relevant *preferences*.

An important question concerns the real world relevance of the findings we have presented in this paper. A key insight we provide is that probability-dependence in risk-taking — if any — will depend on the type of information to which a decision-maker is exposed. Pure descriptions may well result in probability-dependent risk-taking of the standard type. Arguably, however, many real world situations will result in frequent feedback under the form of identically and independently distributed samples, which could induce the opposite type of probability dependence. Many important questions — such as the extent to which some salient events may be more important than less salient events in determining mental probability representations — remain wide open at this point. Finding answers to such questions will prove essential to finding explanations for real-world phenomena, such as lottery play and insurance uptake.

## References

- Abdellaoui, M., L’Haridon, O. & Paraschiv, C. (2011), ‘Experienced versus Described Uncertainty: Do we Need Two Prospect Theory Specifications?’, *Management Science* **57**(10), 1879–1895.
- Atchison, J. & Shen, S. M. (1980), ‘Logistic-normal distributions: Some properties and uses’, *Biometrika* **67**(2), 261–272.
- Aydogan, I. (2021), ‘Prior beliefs and ambiguity attitudes in decision from experience’, *Management Science* **67**(11), 6934–6945.
- Aydogan, I. & Gao, Y. (2020), ‘Experience and rationality under risk: re-examining the impact of sampling experience’, *Experimental economics* **23**(4), 1100–1128.

- Barron, G. & Erev, I. (2003), ‘Small feedback-based decisions and their limited correspondence to description-based decisions’, *Journal of behavioral decision making* **16**(3), 215–233.
- Bishop, C. M. (2006), *Pattern recognition and machine learning*, Vol. 4, Springer.
- Bohren, J. A., Hascher, J., Imas, A., Ungeheuer, M. & Weber, M. (2024), A cognitive foundation for perceiving uncertainty, Technical report, National Bureau of Economic Research.
- Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M., Guo, J., Li, P. & Riddell, A. (2017), ‘Stan: A probabilistic programming language’, *Journal of Statistical Software* **76**(1), 1–32.
- Cubitt, R., Kopsacheilis, O. & Starmer, C. (2022), ‘An inquiry into the nature and causes of the description-experience gap’, *Journal of Risk and Uncertainty* pp. 1–33.
- de Palma, A., Abdellaoui, M., Attanasi, G., Ben-Akiva, M., Erev, I., Fehr-Duda, H., Fok, D., Fox, C. R., Hertwig, R., Picard, N. et al. (2014), ‘Beware of black swans: Taking stock of the description–experience gap in decision under uncertainty’, *Marketing Letters* **25**, 269–280.
- Enke, B. & Graeber, T. (2023), ‘Cognitive uncertainty’, *Quarterly Journal of Economics* **138**(4).
- Erev, I. & Barron, G. (2005), ‘On adaptation, maximization, and reinforcement learning among cognitive strategies.’, *Psychological review* **112**(4), 912.
- Fox, C. R. & Hadar, L. (2006), ‘” decisions from experience” = sampling error+ prospect theory: Reconsidering hertwig, barron, weber & erev (2004)’, *Judgment and Decision Making* **1**(2), 159.
- Frydman, C. & Jin, L. J. (2022), ‘Efficient coding and risky choice’, *Quarterly Journal of Economics* **136**, 161–213.
- Frydman, C. & Jin, L. J. (2025), On the source and instability of probability weighting, Technical report, National Bureau of Economic Research.
- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A. & Rubin, D. B. (2014), *Bayesian data analysis*, Vol. 2, CRC press Boca Raton, FL.



- Glanzer, M., Hilford, A., Kim, K. & Maloney, L. T. (2019), ‘Generality of likelihood ratio decisions’, *Cognition* **191**, 103931.
- Gold, J. I. & Shadlen, M. N. (2001), ‘Neural computations that underlie decisions about sensory stimuli’, *Trends in cognitive sciences* **5**(1), 10–16.
- Gonzalez, R. & Wu, G. (1999), ‘On the Shape of the Probability Weighting Function’, *Cognitive Psychology* **38**(1), 129–166.
- Hau, R., Pleskac, T. J. & Hertwig, R. (2010), ‘Decisions from experience and statistical probabilities: Why they trigger different choices than a priori probabilities’, *Journal of Behavioral Decision Making* **23**(1), 48–68.
- Hau, R., Pleskac, T. J., Kiefer, J. & Hertwig, R. (2008), ‘The description–experience gap in risky choice: The role of sample size and experienced probabilities’, *Journal of Behavioral Decision Making* **21**(5), 493–518.
- Herold, F. & Netzer, N. (2023), ‘Second-best probability weighting’, *Games and Economic Behavior* **138**, 112–125.
- Hertwig, R., Barron, G., Weber, E. U. & Erev, I. (2004), ‘Decisions from experience and the effect of rare events in risky choice’, *Psychological science* **15**(8), 534–539.
- Hertwig, R. & Erev, I. (2009), ‘The description–experience gap in risky choice’, *Trends in cognitive sciences* **13**(12), 517–523.
- Hertwig, R. & Pleskac, T. J. (2010), ‘Decisions from experience: Why small samples?’, *Cognition* **115**(2), 225–237.
- Imai, T., Nunnari, S., Wu, J. & Vieider, F. M. (2025), ‘Meta-analysis of prospect theory parameters’, *Working Paper*.
- Jessup, R. K., Bishara, A. J. & Busemeyer, J. R. (2008), ‘Feedback produces divergence from prospect theory in descriptive choice’, *Psychological Science* **19**(10), 1015–1022.
- Kahneman, D. & Tversky, A. (1979), ‘Prospect Theory: An Analysis of Decision under Risk’, *Econometrica* **47**(2), 263 – 291.
- Khaw, M. W., Li, Z. & Woodford, M. (2021), ‘Cognitive imprecision and small-stakes risk aversion’, *The Review of Economic Studies* **88**(4), 1979–2013.

- Khaw, M. W., Li, Z. & Woodford, M. (2025), Cognitive imprecision and stake-dependent risk attitudes, Technical report.
- Ma, W. J., Kording, K. P. & Goldreich, D. (2023), *Bayesian Models of Perception and Action: An Introduction*, MIT press.
- Natenzon, P. (2019), ‘Random choice and learning’, *Journal of Political Economy* **127**(1), 419–457.
- Netzer, N. (2009), ‘Evolution of time preferences and attitudes toward risk’, *American Economic Review* **99**(3), 937–55.
- Netzer, N., Robson, A., Steiner, J. & Kocourek, P. (2024), ‘Risk perception: Measurement and aggregation’, *Journal of the European Economic Association* **In Press**.
- Olschewski, S. & Scheibehenne, B. (2024), ‘What’s in a sample? epistemic uncertainty and metacognitive awareness in risk taking’, *Cognitive Psychology* **149**, 101642.
- Oprea, R. (2024), ‘Decisions under risk are decisions under complexity’, *American Economic Review* **114**, 3789–3811.
- Preston, M. G. & Baratta, P. (1948), ‘An Experimental Study of the Auction-Value of an Uncertain Outcome’, *The American Journal of Psychology* **61**(2), 183.
- Robson, A. J. (2001*a*), ‘The biological basis of economic behavior’, *Journal of Economic Literature* **39**(1), 11–33.
- Robson, A. J. (2001*b*), ‘Why would nature give individuals utility functions?’, *Journal of Political Economy* **109**(4), 900–914.
- Robson, A. J. & Samuelson, L. (2011), The evolutionary foundations of preferences, in ‘Handbook of social economics’, Vol. 1, Elsevier, pp. 221–310.
- Savage, L. J. (1954), *The Foundations of Statistics*, Wiley, New York.
- Simon, H. A. (1959), ‘Theories of decision-making in economics and behavioral science’, *The American Economic Review* **49**(3), 253–283.
- Steiner, J. & Stewart, C. (2016), ‘Perceiving prospects properly’, *American Economic Review* **106**(7), 1601–31.

- Tversky, A. & Kahneman, D. (1992), ‘Advances in Prospect Theory: Cumulative Representation of Uncertainty’, *Journal of Risk and Uncertainty* **5**, 297–323.
- Tversky, A. & Wakker, P. P. (1995), ‘Risk Attitudes and Decision Weights’, *Econometrica* **63**(6), 1255–1280.
- Tymula, A., Wang, X., Imaizumi, Y., Kawai, T., Kunimatsu, J., Matsumoto, M. & Yamada, H. (2023), ‘Dynamic prospect theory: Two core decision theories coexist in the gambling behavior of monkeys and humans’, *Science Advances* **9**(20), eade7972.
- Ungemach, C., Chater, N. & Stewart, N. (2009), ‘Are probabilities overweighted or underweighted when rare outcomes are experienced (rarely)?’, *Psychological Science* **20**(4), 473–479.
- Van de Kuilen, G. (2009), ‘Subjective probability weighting and the discovered preference hypothesis’, *Theory and decision* **67**, 1–22.
- van de Kuilen, G. v. d. & Wakker, P. P. (2006), ‘Learning in the allais paradox’, *Journal of Risk and Uncertainty* **33**(3), 155–164.
- Vieider, F. M. (2024a), Bayesian estimation of decision models, Technical report.  
**URL:** <https://fvieider.quarto.pub/bstats/>
- Vieider, F. M. (2024b), ‘Decisions under uncertainty as bayesian inference on choice options’, *Management Science* **70**, 8217–9119.
- von Neumann, J. & Morgenstern, O. (1944), *Theory of Games and Economic Behavior*, Princeton University Press, New Heaven.
- Wakker, P. P. (2010), *Prospect Theory for Risk and Ambiguity*, Cambridge University Press, Cambridge.
- Wulff, D. U., Mergenthaler-Canseco, M. & Hertwig, R. (2018), ‘A meta-analytic review of two modes of learning and the description-experience gap.’, *Psychological bulletin* **144**(2), 140.
- Zhang, H. & Maloney, L. T. (2012), ‘Ubiquitous log odds: a common representation of probability and frequency distortion in perception, action, and cognition’, *Frontiers in neuroscience* **6**, 1.
- Zhang, H., Ren, X. & Maloney, L. T. (2020), ‘The bounded rationality of probability distortion’, *Proceedings of the National Academy of Sciences* **117**(36), 22024–22034.

## Online Appendices

## A Model derivation

### A.1 A general inference model

In Decisions from Experience (DfE), subjects need not only learn the outcomes and underlying probabilities, but also the whole structure of the decision problem (i.e., the number of outcomes in the lottery's support). In the body of the paper we assume away this component of the inference problem for simplicity and to focus our discussion on the influence of sampling error and sampling variance. Here, for completeness, we propose a stylized model of how such higher order learning could take place based on the sort of sampling from the two options that occurs in DfE. We argue that expanding the model in this way has little qualitative impact on our findings.

We start by discussing the structural inference process. Assume a DM believes that outcomes will range from 0 to some upper limit  $u$ , outcomes beyond which are not considered plausible.<sup>27</sup> Take two objective probability distributions over all outcomes underlying the two choice options,  $\{p_0, p_1, \dots, p_u\}$  and  $\{q_0, q_1, \dots, q_u\}$ , where subscripts indicate monetary outcomes. In DfE, DMs will infer the probability distributions from the draws they observe. Let the initial likelihood at time  $t = 0$ , before any draws are taken, be encoded in two  $u + 1$ -dimensional Dirichlet distributions,  $\mathcal{D}_A(\pi_j) \propto \prod_{j=0}^u \tilde{p}_j^{\pi_j - 1}$  and  $\mathcal{D}_B(\omega_j) \propto \prod_{j=0}^u \tilde{q}_j^{\omega_j - 1}$ , where  $\tilde{p}_i \triangleq \frac{\pi_i}{\sum_j \pi_j}$  and  $\tilde{q}_i \triangleq \frac{\omega_i}{\sum_j \omega_j}$  represent the subjective expectations of the probabilities attributed to an outcome  $i$  in the two choice options  $A$  and  $B$ . Given the ex ante exchangeability of the two choice options, the two Dirichlets will have the same parameters at time  $t = 0$ . We assume that DMs consider any given outcome as equally likely in the two choice options, so that  $\pi_i = \omega_i \forall i$  at  $t = 0$ . This assumption directly follows from the exchangeability of the two options before any draws have been observed, and is implemented in our experiment by randomizing the risky and safe options in positions A and B.

We assume that what matters for decisions is the direct comparison between the two choice options. To capture this in our model, we map the inferences based on the Dirichlets encoding draws from the two choice options into a *comparative Dirichlet* which entails a statewise comparison between two options. That is, what matters for choices are events in which one option pays a given outcome, while the other option pays a different outcome. In our experiment, these will be the events under which the risky option pays  $x$  while the

---

<sup>27</sup>In principle,  $u$  can take any value, as long as it is finite. In our experiment, we tell subjects beforehand that all outcomes will range between \$0 and \$35 inclusive, thus setting their expectations about this range.

safe option pays  $c < x$ , and the event under which the risky option pays  $y$  while the safe option pays  $c > y$  (see below for a generalization). The probabilities of the comparative events  $e_1$  (obtain  $x > c$  rather than  $c$ ) and  $e_2$  (obtain  $c$  rather than  $y < c$ ) can now be obtained from the single-state Dirichlets  $\mathcal{D}_A(\pi_j)$  and  $\mathcal{D}_B(\omega_j)$  defined for the two options, since  $P[e_1] = P[x \cap c] = \tilde{p}_x \times \tilde{p}_c$  and  $P[e_2] = P[y \cap c] = \tilde{p}_y \times \tilde{p}_c$ . Given that for finite samples  $\tilde{p}_c < 1$  and  $\tilde{p}_x + \tilde{p}_y < 1$ , the inferred probabilities will generally be subadditive, that is,  $P[e_1] + P[e_2] \leq 1$  (with 1 being the limiting case as samples tend to infinity). This implies that we can express the subjective beliefs in the comparative states of the world once again by a Dirichlet,  $\mathcal{D}(\delta_i) = \prod_{i=1}^u P[e_i]^{\lambda \hat{\delta}_i - 1}$ , where  $\lambda \triangleq \sum_{i=1}^u \delta_i$  is the concentration of the new Dirichlet, and  $\hat{\delta}_i \triangleq \delta_i / \lambda$  captures the mean belief about a given state  $i$ . While some probability mass will thus remain attributed to ‘non-observed outcomes’, this part will drop out of the main choice equation below.

This justifies the assumption of the Beta distribution in the main text: while the latter imposes additivity in  $\hat{p}_x$  and  $\hat{p}_y$ , that assumption serves to simplify our discussion, but has no substantive implications for our conclusions (given that the non-observed states receiving the remaining probability mass drop out of the discriminability equation). If, say, a third outcome from the risky option were to be observed at some point, this would add a new comparative state to the comparison (see below). In the text we further discussed inference bias in terms of the samples taken from the risky option only. More generally, however, the samples from the safe option will also count. While a precise closed-form solution does not exist for that case, we can approximate the samples by the total samples for each state, where the samples from the safe option are simply added to the samples indicating each comparative sample in the sum of the trigamma functions. This means that our discussion in the main text may *quantitatively* underestimate the samples, but that this more general case will not qualitatively affect any of the conclusions drawn.

In the main text, we implicitly assume that subjects know which of the two options is the risky one and which the safe. In reality, subjects need to infer this from the samples they take. We make three assumptions in this regard. The first, and most substantively relevant, is that subjects make inferences on the choice environment (including potentially the intentions of the experimenter). This entails that choices between two non-degenerate options are deemed extremely unlikely. Practically, this entails that sampling variance will remain high until a plausible set of outcomes has been observed.<sup>28</sup> The second assumption

---

<sup>28</sup>This assumption seems particularly defensible in our DfE experiments, since all subjects assigned to this treatment have all finished making dozens of binary DfD choices for lotteries with one degenerate and one

is that we assume the initial parameters of the two choice option Dirichlets to be sparse, i.e.  $\pi_i, \omega_i \ll 1 \forall i$ . This assumption implies that subjects do not expect a very diffuse probability distribution with many different outcomes. Practically, this helps explain why samples are relatively small, since it keeps the probability mass assigned to unobserved outcomes low in the comparative Dirichlet.

An additional assumption in the main text is that subjects can infer which of the two options is the risky one. This obtains trivially once a subject has observed all three outcomes used in our experiment (the two in the risky option, and the one in the safe option, which constitute a ‘plausible minimal outcome set’ inasmuch as they indicate a non-degenerate choice, or equivalently, they map into two comparative states with a meaningful tradeoff between log-odds and log-cost benefits). This indeed follows directly from the two assumptions above: that subjects expect non-degenerate choices, and that the initial parameters are sparse (meaning that they do not necessarily expect more outcomes once they have observed a plausible outcome set). The inference is somewhat less trivial as long as only one outcome has been observed from each choice option.

We illustrate this based on the choice options we provide in the experiment. For small probabilities, subjects are overwhelmingly likely to observe the lower outcome  $y$ . Given that in our experiment  $y$  is always equal to 0, and that we tell subjects that they will only ever face non-negative amounts, this immediately identifies this choice option as the risky one. For large probabilities, where subjects may observe two strictly positive amounts  $x$  and  $c$  from the two options, this is less obvious. We thus furthermore assume that the parameters of the option-specific Dirichlets before any samples are taken will be characterized by sparsity increasing in outcomes. That is, for any  $j > i$ , where the two indices are non-negative outcomes,  $\omega_j = \pi_j \leq \pi_i = \omega_i$  at time  $t = 0$ , before any samples have been taken. In practice, this entails that subjects consider smaller outcomes more likely than larger outcomes. Notice that this is the equivalent of a pessimistic prior for the inference process, and that it is thus fully coherent with both our model and our empirical results.<sup>29</sup>

---

non-degenerate lottery.

<sup>29</sup>In principle, this inference process could be modelled as a probabilistic process resulting in stochastic assessments of the riskiness of the two choice options after each sample. Such a model would follow a very similar structure as our discriminability model, and we do thus not formalize it here. Such a model would be most relevant for large probability lotteries in cases where only one outcome has been observed from each option. The notion that subjects infer the structure of such choice problems from sampling draws is indeed supported by the observation that samples from the *safe* option increase in the objective probability of winning for both risk averse and risk seeking subjects in our data.

## A.2 Noisy log-odds representation

In our actual experiment, subjects will experience exactly 1 outcome from the sure option, and no more than 2 from the risky option. We can thus use the 2-dimensional special case of the comparative Dirichlet distribution discussed above – the Beta distribution (see above for an explicit discussion of this simplifying assumption). In particular, the parameter  $\alpha$  will encode the ‘good state’, in which the lottery pays a prize  $x > c$ , whereas  $\beta$  will encode the ‘bad state’, under which the lottery pays an outcome  $y < c$ . The perceived or sampled probability of the good state favoring the lottery will thus be  $\mathbb{E}[\hat{p}] = \frac{\alpha}{\alpha+\beta}$ .

We start from an optimal choice rule entailing expected value maximization. The DM will thus choose the lottery over the sure amount whenever  $\hat{p}x + (1 - \hat{p})y > c$ , or equivalently whenever

$$\ln \left( \frac{\hat{p}}{1 - \hat{p}} \right) > \ln \left( \frac{c - y}{x - c} \right).$$

The transformation into log-odd space is convenient for computational reasons, but otherwise inconsequential (see Vieider 2024b, for an alternative derivation). The choice rule entails that the log-odds in favor of the lottery will be traded off against the log of the ratio of costs ( $c - y$ , potentially get the lower outcome  $y$  when  $c$  could have been had) and benefits ( $x - c$ ; obtain the prize  $x$  instead of the lower sure amount  $c$ ). Here, we will assume without loss of generality that the log cost-benefits are perceived objectively. This is a simplifying assumption that allows us to focus on the likelihood dimension, where most of the action takes place. It is straightforward to generalize the derivation to include the noisy coding of costs and benefits as well (cfr. Vieider 2024b).

The mean of the sampled log-odds can simply be derived from the two parameters containing the counts of successes and failures:

$$\mathbb{E} \left[ \ln \left( \frac{\hat{p}}{1 - \hat{p}} \right) \right] = \ln \left( \frac{\alpha}{\beta} \right)$$

Given limited samples, however, even samples that are accurate on average will contain some error on single draws, driven by natural sampling variation around the true mean. Averaging across all probabilities, we will thus observe

$$\ln \left( \frac{\alpha}{\beta} \right) = \ln \left( \frac{p}{1 - p} \right) + \varepsilon,$$

which, following Atchison & Shen (1980), could equivalently be written as the difference of



the digamma functions of the two parameters,  $F(\alpha) - F(\beta)$ .

Log-odds tend to follow approximately normal distributions, giving rise to a *logit-normal* (Atchison & Shen 1980). This suggests that  $\varepsilon \sim \mathcal{N}(0, \nu^2)$ . The sampling variance  $\nu^2$ , in turn, again derives from the properties of the logit-normal distribution, and is given by the sum of trigamma functions of the two parameters, i.e.  $\nu^2 = F'(\alpha) + F'(\beta)$ .

## Optimal Combination with Bayesian prior

Given the noise in inferences, it will be optimal to combine the observations with a Bayesian prior. The optimality of this operation derives from the fact that — even though it will introduce systematic bias into the estimates under the form of regression to the mean of the prior — it will minimize the mean squared error across many estimates (see Ma et al. 2023, chapter 4, for an illustration). The reason for this is that the reduction in variance of the estimator will more than make up for the introduction of bias.

The objective for the mind now becomes to infer the log-odds from the underlying samples (whether they be true samples or virtual/neural samples — we drop the subscripts here and derive the equation just once). The inference problem for any given choice task will thus be as follows:

$$\mathbb{E} \left[ \ln \left( \frac{p}{1-p} \right) \mid \alpha, \beta \right] = \frac{\sigma^2}{\sigma^2 + \nu^2} \ln \left( \frac{\alpha}{\beta} \right) + \frac{\nu^2}{\sigma^2 + \nu^2} \mu,$$

where we redefine  $\mu = \ln \left( \frac{p_0}{1-p_0} \right)$  in the main text, and where the Bayesian evidence weight or “likelihood-discriminability” parameter is given by  $\gamma \triangleq \frac{\sigma^2}{\sigma^2 + \nu^2} = \frac{1}{1 + \nu^2/\sigma^2}$ . A step-by-step derivation of this equation can be found in Vieider (2024a), chapter 2.

In DfD, the “virtual draws” encoded in  $\alpha$  and  $\beta$  (referred to as  $\hat{\alpha}$  and  $\hat{\beta}$  in the main text) are unobservable. We can, however, estimate the equation by aggregating across multiple similar probabilities. This will yield the expectation over repeated stimuli of the posterior expectation above, which takes the following form:

$$\mathbb{E} \left[ \mathbb{E} \left[ \ln \left( \frac{p}{1-p} \right) \mid \frac{\hat{\alpha}}{\hat{\beta}} \right] \mid p \right] = \frac{\sigma^2}{\sigma^2 + \nu^2} \ln \left( \frac{p}{1-p} \right) + \frac{\nu^2}{\sigma^2 + \nu^2} \mu,$$

which now allows us to substitute the true log-odds for the sampled log-odds. Choice to choice fluctuations in the samples will be reflected in the variance of the distribution, which

takes the form  $\gamma^2 \nu^2 = \frac{\sigma^4 \nu^2}{(\sigma^2 + \nu^2)^2}$ .

*Proof.* The proof exploits the well-known property of the normal distribution whereby  $z \sim \mathcal{N}(\hat{z}, \tau^2)$  implies  $bz + a \sim \mathcal{N}(b\hat{z} + a, b^2\tau^2)$ . To obtain the response distribution above, let  $\ln\left(\frac{\alpha}{\beta}\right) = z$ ,  $\frac{\sigma^2}{\sigma^2 + \nu^2} = b$ ,  $\frac{\nu^2}{\sigma^2 + \nu^2} \mu = b$ ,  $\ln\left(\frac{p}{1-p}\right) = \hat{z}$ , and  $\nu = \tau$ .  $\square$

Note that the problem does not change in any substantive way if we abandon the assumption of draws correctly reflecting the underlying distribution on average when real samples are taken in DfE. We then simply change the objective probability  $p$  to the sampled probability  $\hat{p}$  in the equations above. Sampling bias in  $\hat{p}$  will then occur on top of the inference bias, which still results in regression to the mean of the prior, just like represented above.

## Stochastic choice rule

We can now trade off the inferred log-odds, as derived above, against the log-cost benefits, as suggested by our optimal choice rule. Letting  $\mu \triangleq \ln\left(\frac{p_0}{1-p_0}\right)$ , we obtain  $\delta = \ln\left(\frac{p_0}{1-p_0}\right)^{1-\gamma}$ , and by extension,  $\theta = \delta^{-1} = \ln\left(\frac{1-p_0}{p_0}\right)^{1-\gamma}$ . Putting everything on the scale of the standard deviation of the response distribution derived in the previous section yields the z-score describing the choice probability of the lottery:

$$pr[(x, p; y) \succ c] = \Phi \left[ \frac{\gamma \ln\left(\frac{p}{1-p}\right) - \ln\left(\frac{c-y}{x-c}\right) - \ln(\theta)}{\gamma \nu} \right],$$

where  $\Phi$  is the standard normal cumulative distribution function. In DfD (as well as DfD+forced and DfE+forced), the probability will correspond to the correct one, and the model can thus be simply estimated on choice data by plugging the probit link function above into a Bernuoulli distribution (see below).

In DfE, we need to slightly amend the function above. In particular, we will now substitute sampled probabilities  $\hat{p}$  for the true probabilities above (adding a constant to both numerator and denominator to make sure it is defined—see discussion of the inference process above). An additional assumption concerns the log cost-benefit ratio when either  $x$ ,  $y$ , or  $c$  have not yet been observed. The simplest assumption is that of a “naive” decision maker, who assumes the ratio to be 1 in that case (and hence its logarithm to be 0). However, this is just a special case of what a more sophisticated decision maker would do. Multiplying the log cost-benefit ratio by an additional parameter  $\rho$ , conditional on one of the outcomes

not yet having been observed, allows for a more flexible specification whereby the DMs can (correctly) infer a positive correlation between log-odds and log cost-benefits. The “naive” DM discussed above is then just a special case for whom  $\rho = 0$ .

## N-dimensional generalization

The inference framework discussed at the beginning of this section is fully general. While we have described it for the particular case of comparisons used in our experiment, it can just as easily be applied to comparison between multi-outcome lotteries. The inference framework introduced above remains directly applicable, with the two option-specific Dirichlet simply counting instances of different outcomes. Our setup assumes that outcomes are ordered by size to arrive at the comparative distribution. The comparative Dirichlet is then constructed over  $k$  comparative states constructed based on the ranked outcomes.

Take two lotteries offering outcomes  $\mathbf{x} = \{x_1, \dots, x_k\}$  and  $\mathbf{y} = \{y_1, \dots, y_k\}$  under the comparative events  $e_1, \dots, e_k$ , where each comparative event is characterized by a probability  $\hat{p}_i$ , which could be different from the true underlying probability  $p_i$ . We assume that the outcome are ordered such that  $x_1 \geq x_2 \geq \dots \geq x_k$  and  $y_1 \geq y_2 \geq \dots \geq y_k$ . We further assume for our representation that  $\mathbf{x}$  is riskier than  $\mathbf{y}$  in the sense of having wider spread or variance. Draws from the two choice options ought to be seen as independent, just as is the case in the actual samples taken. The optimal choice rule, which once again entails expected value maximization, takes the following form:

$$\sum_{i=1}^k \frac{\hat{p}_i}{1 - \hat{p}_i} (x_i - y_i) > 1, \quad (5)$$

which sums the relative benefits of the riskier option,  $x_i - y_i$ .

Assuming that the different states will be processed in parallel, the stochastic choice equation then takes the following form:

$$P[\mathbf{x} \succ \mathbf{y}] = \sum_{i=1}^k \Phi \left[ \frac{\gamma \times \ln \left( \frac{\hat{p}_i}{1 - \hat{p}_i} \right) + \mathbb{1} \times \ln (\mathbb{1}(x_i - y_i)) - \ln(\theta)}{(k-1) \nu \times \gamma} \right].$$

where  $\mathbb{1} = 1$  if  $x_i - y_i > 0$  and else  $\mathbb{1} = -1$ , thus assuring that the logarithm is defined. The multiplication of the “relative benefit” by  $\mathbb{1}$  further makes sure that this quantity enters with the appropriate sign, since it could favor either choice option in any given state  $i$ . Given

that any single comparison is standard-normally distributed, the sum over the different comparisons will also follow a standard normal distribution. While this formulation *could*, in principle, result in predicted choice probabilities greater than 1 or smaller than 0, this is unlikely in practice, given that benefits and costs are usually designed to compensate each other. A regularizing condition could be imposed to overcome this issue should it ever become relevant in practice. We have careful study of this extension for future work.

## B Experiments

### Choice stimuli

We selected our choice stimuli from those in the early DfE literature (Hertwig et al. 2004), but generalized them so as to allow us to structurally estimate our model, and to obtain a more balanced picture of the behavior. We assured identification of the structural estimations using simulations, which allowed us to find the optimal compromise between number and type of task and the length of the experiment. The limiting factor derived in particular from the forced sampling experiments, where subjects had to take 40 samples by tasks, as well as expressing their final choice.

We thus chose 6 different lotteries—3 with a small probability, and 3 with a large probability of winning. We then obtained three choice tasks by lottery by setting the sure amount  $c$  equal to the expected value, and by adding or subtracting a fixed amount. This provides some valuable variation for the structural estimations, and results in the following 18 unique tasks (4 randomly selected ones of which were repeated in the experiment):

## C Meta-analytic estimation

**Quantifying the GAP.** To get a better idea of the size of the decision-experience GAP in our data, and to relate it to typical findings in the literature, we can aggregate the evidence across tasks using the tools of meta-analysis.<sup>30</sup> Let  $\pi_d = R_d/N_d$  be the proportion of risky choices in DfD, where  $R_d$  is the number of risky choices, and  $N_d$  the number of observations.

---

<sup>30</sup>The meta-analytic tools we use are identical to a “measurement error model”. That is, the assumption is that each single choice proportion is observed with some error. Meta-analysis then allows us to aggregate across the choice proportions while eliminating measurement error and thus correcting our analysis for multiple testing across many moderate (and not statistically independent) samples.

Table 3: Choice tasks

| small $p$         | large $p$         |
|-------------------|-------------------|
| (31,0.10) vs. 2.8 | (4,0.80) vs. 3.0  |
| (31,0.10) vs. 3.2 | (4,0.80) vs. 3.2  |
| (31,0.10) vs. 3.6 | (4,0.80) vs. 3.4  |
| (10,0.15) vs. 1.2 | (8,0.85) vs. 6.6  |
| (10,0.15) vs. 1.5 | (8,0.85) vs. 6.8  |
| (10,0.15) vs. 1.8 | (8,0.85) vs. 7.0  |
| (16,0.20) vs. 2.9 | (10,0.90) vs. 8.8 |
| (16,0.20) vs. 3.2 | (10,0.90) vs. 9.0 |
| (16,0.20) vs. 3.5 | (10,0.90) vs. 9.2 |

Choice tasks are describes as usual, with  $(x, p)$  designating a lottery providing a prize  $x$  with probability  $p$  or else 0, and  $c$  designating the sure amount.

Let  $\pi_e = R_e/N_e$  be the proportion in DfE. We define the difference in choice proportions as  $g$ , where we encode the difference in the direction of the standard gap, so that  $g = \pi_d - \pi_e$  for  $p < 0.5$  and  $g = \pi_e - \pi_d$  for  $p > 0.5$ . This difference will be approximately normally distributed, with variance  $\pi(1 - \pi)(1/N_d + 1/N_e)$ , where  $\pi = \frac{\pi_d + \pi_e}{N_d + N_e}$ . We can now use  $g$  and its associated standard error,  $se$ , for meta-analytic aggregation across tasks, indexed by  $i$ :

$$g_i \sim \mathcal{N}(\hat{g}_i, se_i^2)$$

$$\hat{g}_i \sim \mathcal{N}(\omega, \tau^2),$$

where  $g$  and  $se$  are data,  $\hat{g}$  is the unknown true effect, and  $\omega$  and  $\tau$  are parameters capturing the meta-analytic mean and standard deviation across tasks, respectively. We then quantify the GAP by meta-analytically aggregating the differences in choice proportions across tasks in a direction that is consistent with the standard GAP.

**Reversals in Likelihood Dependence.** We can also use meta-analysis to test whether choice proportions exhibit probability-dependence, and whether the nature of this dependence is different in DfE and DfD. To do this, we analyze the choice proportions  $\pi_i$  directly (instead of examining differences in choice proportions  $g_i$ ) so that we estimate  $\pi_i \sim \mathcal{N}(\hat{\pi}_i, se_i)$ . We then use meta-regression to assess the dependence of the choice proportions on the probability of winning, by letting  $\hat{\pi}_i \sim \mathcal{N}(\lambda_0 + \lambda \times p_i, \tau^2)$ , where  $\hat{\pi}_i$  is the unknown true choice proportion.

We estimate the model in Stan (see Vieider 2024a for a tutorial on the use of Stan for decision models; chapter 4 contains a part specifically dedicated to meta-analysis). Here is

the Stan code used to estimate the model:

```
//footnotesize
data{
  int<lower=1> N; \\number of observation
  vector[N] gap; \\difference in choice proportions
  vector<lower=0>[N] se; \\standard error of the difference
}
parameters{
  vector[N] gamma; //true, estimated gap (called g_hat in paper)
  real mu; //meta-analytic mean (omega in paper)
  real<lower=0> sigma; //variance
}
model{
  //regularizing priors
  sigma ~ normal( 0 , 1 );
  mu ~ normal( 0 , 1 );

  // measurement error model:
  gap ~ normal( gamma , se );

  // likelihood:
  gamma ~ normal( mu , sigma );
}
```

The meta-regression is introduced into the same code simply by modifying the mean  $\mu$ , making it dependent on the probability of winning:

```
//footnotesize
data{
  int<lower=1> N;
  int<lower=1> K; //dimension of design matrix
  vector[N] gap;
  vector<lower=0>[N] se;
  matrix[N,K] X; //design matrix of explanatory variables
}
```

```

parameters{
  vector[N] gamma;
  real mu;
  real<lower=0> sigma;
  vector[K] beta;
}
model{
  sigma ~ normal( 0 , 1 );
  mu ~ normal( 0 , 1 );

  // measurement error model:
  gap ~ normal( gamma , se );
  // likelihood:
  gamma ~ normal( mu + X * beta , sigma );
}

```

## D Structural estimation

We implement our structural equations based on the discriminability equation in the main text, using the objective probability of winning,  $p$ , in DfD, DfD+forced, and DfE+forced. We use the sampled probability  $\ln\left(\frac{\alpha}{\beta}\right)$  in DfE, and complement this with an assumption about the log-cost benefits in the case that one of the outcomes has not yet been observed when the decision is taken, as described above.

We keep the model as simple as possible in order to maximize our comparative power and to keep the model parsimonious. This means, first of all, that we normalize the coding noise variance by division with the variance of the prior, so that  $\gamma = \frac{1}{1 + \frac{\nu^2}{\sigma^2}}$ . This helps both identifiability and comparability across treatments but happens without loss of generality, since it is the ratio between coding noise variance and prior variance that determines behavior (see also Natenzon 2019). Another assumption that we maintain throughout the paper is that the mean of the prior,  $\mu$ , remains unaffected over the course of the experiment. We exploit this in the estimation by letting  $\mu$  be the same across the 2 parts of the experiment,

whereas  $\nu$  and as a consequence  $\gamma$  and  $\theta$  are all allowed to vary freely.

We estimate the model using a Bayesian hierarchical setting in Stan (Carpenter et al. 2017). The hierarchical setting allows us to pool information from the aggregate estimation, which provides the priors, and from individual-level parameter estimates, which contribute to the aggregate in proportion to their precision. The aggregation equation follows exactly the equation we describe for our Bayesian inference process. Vieider (2024a) provides a step-by-step tutorial on the estimation of decision models in Stan.

Below, we include an commented version of the code we use in DfD, DfD+forced, and DfE+forced (the code used in DfE is very similar, and only has an additional parameter  $\rho$ , as well as including the truly observed log-odds as data; it is available upon request). We define the variables at the level of the individual *choices*. This allows us to implement a literal specification of our model, where task-specific quantities are encoded by parameters  $\alpha$  and  $\beta$ . These parameters are nested in individual-level parameters, which we use to fit the choice data, and which ensures that the choice-level parameters are identified and well-behaved (since the individual-level parameters act as informative priors). Finally, individual-level parameters are nested within an overall distribution.

We check convergence by making sure that all R-hats are below 1.05. We also carefully check that any divergent iterations do not indicate problems with the posterior (and discard all estimates with more than 1% divergent iterations). The hyperpriors on the aggregate parameter means are given very wide priors, which makes them *mildly regularizing*—they help the convergence of the simulation algorithm by being centered around the region where we expect the parameter values to fall, but they attribute significant probability mass to 1 order of magnitude above the region into which we would expect the parameters to reasonably fall. Our estimates are indeed not sensitive to the choice of the exact parameter values. This follows best practices in Bayesian estimation.

```
data{
  \\\declare data
  int<lower=1> N; \\\number of observations
  int<lower=1> N_id; \\\number of subjects
  array[N] int id; \\\unique identifier
  array[N] real high; \\\outcome x
  array[N] real low; \\\outcome y
  array[N] real sure; \\\outcome c
  array[N] real p; \\\probability
  array[N] int choice_risky; \\\choice: 1 if risky
  array[N] int part2; \\\dummy to indicate part 2
```



```

}
transformed data{
  array[N] real lcb; \\log cost benefit ratio
  array[N] real llr; \\log-odds
  for (i in 1:N){
    lcb[i] = log( (sure[i] - low[i]) / (high[i] - sure[i]) );
    llr[i] = log( p[i]/(1 - p[i]) );
  }
}
parameters{
  vector[3] means; \\aggregate mean parameters on log scale
  vector<lower=0>[3] tau_id; \\aggregate parameter variances
  cholesky_factor_corr[3] L_omega_id; \\decomposed covar matrix
  array[N_id] vector[3] Zid; \\stan dardized individual-level parameters
}
transformed parameters{
  // covar and temp parameters
  matrix[3,3] Rho_id = L_omega_id * L_omega_id'; \\obtain covariance matrix
  array[N] vector[3] pars; \\parameter matrix on log scale
  // generative parameters:
  vector[N] mu; \\prior mean
  vector<lower=0>[N] kappa; \\concentration part1
  vector<lower=0>[N] kappaf; \\concentration part2
  // derived parameters from here
  vector[N] alpha; \\derived parameters—see definitions in text, and below
  vector[N] beta;
  vector[N] nu;
  vector[N] gamma;
  vector[N] theta;
  vector[N] omega;
  vector[N] alphaf;
  vector[N] betaf;
  vector[N] nuf;
  vector[N] gammaf;
  vector[N] thetaf;
  vector[N] omegaf;
  for (i in 1:N){
    pars[i] = means + diag_pre_multiply(tau_id, L_omega_id) * Zid[id[i]];
    mu[i] = pars[i,1];
    kappa[i] = exp(pars[i,2]);
    kappaf[i] = exp(pars[i,3]);
  }
  // define derived parameters
  alpha[i] = kappa[i] * p[i];

```

```

    beta[i] = kappa[i] * (1 - p[i]);
    nu[i] = sqrt( trigamma( alpha[i] ) + trigamma( beta[i] ) );
    gamma[i] = 1/( 1 + nu[i]^2 );
    theta[i] = exp( ( gamma[i] - 1 ) * mu[i] ) ;
    omega[i] = nu[i] * gamma[i];
    alphaf[i] = kappaf[i] * p[i];
    betaf[i] = kappaf[i] * (1 - p[i]);
    nuf[i] = sqrt( trigamma( alphaf[i] ) + trigamma( betaf[i] ) );
    gammaf[i] = 1/( 1 + nuf[i]^2 );
    thetاف[i] = exp( ( gammaf[i] - 1 ) * mu[i] ) ;
    omegaf[i] = nuf[i] * gammaf[i];
  }
}
model{
  vector[N] udiff; \\local vector
  \\priors for aggregate (hierarchical) parameters
  tau_id ~ exponential(5);
  L_omega_id ~ lkj_corr_cholesky(4);
  means[1] ~ normal(0, 5);
  means[2] ~ normal(0, 5);
  means[3] ~ normal(0, 5);

  \\priors for individual level parameters, standardized:
  for (n in 1:N_id)
    Zid[n] ~ std_normal();

  \\the mode:
  for ( i in 1:N ) {
    udiff[i] = ( ( gamma[i] * llr[i] - lcb[i] - log(theta[i]) )/ omega[i] ) * (1 - part2[i])
               + ( ( gammaf[i] * llr[i] - lcb[i] - log(thetaf[i]) )/ omegaf[i] ) * part2[i] );
    choice_risky[i] ~ bernoulli( Phi( udiff[i] ) );
  }
}
\\code below recovers individual-level parameters
generated quantities{
  vector[N] log_lik;
  vector[N] udiff;

  vector[N_id] mun;
  vector[N_id] kappan;
  vector[N_id] alphan;
  vector[N_id] betan;
  vector[N_id] nun;

```

```

vector[N_id] gamman;
vector[N_id] thetan;
  vector[N_id] kappafn;
vector[N_id] alphafn;
vector[N_id] betafn;
vector[N_id] nufn;
vector[N_id] gammafn;
vector[N_id] thetafn;

vector[3] temp;
for(n in 1:N_id){
  temp = means + diag_pre_multiply(tau_id, L_omega_id) * Zid[n];
  mun[n] = temp[1];
  kappan[n] = exp(temp[2]);
  kappafn[n] = exp(temp[3]);
  alphan[n] = kappan[n]/2;
  betan[n] = kappan[n]/2;
  nun[n] = sqrt( trigamma( alphan[n] ) + trigamma( betan[n] ) );
  gamman[n] = 1/(1 + nun[n]^2 );
  thetan[n] = exp( ( gamman[n] - 1 ) * mun[n] );
  alphafn[n] = kappafn[n]/2;
  betafn[n] = kappafn[n]/2;
  nufn[n] = sqrt( trigamma( alphafn[n] ) + trigamma( betafn[n] ) );
  gammafn[n] = 1/(1 + nufn[n]^2 );
  thetafn[n] = exp( ( gammafn[n] - 1 ) * mun[n] );
}

for ( i in 1:N ) {
  udiff[i] = ( ( gamma[i] * llr[i] - lcb[i] - log(theta[i]) ) / omega[i] ) * (1 - comp[i])
              ( ( gammaf[i] * llr[i] - lcb[i] - log(thetaf[i]) ) / omegaf[i] ) * comp[i];
  log_lik[i] = bernoulli_lpmf( choice_risky[i] | Phi_approx( udiff[i] ) );
}
}

```

## D.1 Structural estimation results

We use structural estimation to more deeply assess the hypothesis that both probability-dependence and the description-experience gap are a consequence of cognitive noise — and that our treatments eliminate these patterns by eliminating this noise. We structurally estimate our model from choice data based on our discriminability equation (4). The key

parameter driving both probability-dependence and the GAP in our model (and, therefore, our focus in this section) is  $\gamma$ , the weight the DM puts on her perception of the log-odds in the decision process. We will refer to this as “likelihood-discriminability,” mirroring the name given the equivalent parameter in the LLO function, “likelihood-sensitivity.” In the model,  $\gamma$  is an inverse function of coding noise: the smaller coding noise  $\nu$  becomes, the closer  $\gamma$  will come to 1, producing perfect discriminability of log-odds. Importantly, this parameter is estimated, in part, using inconsistencies in subjects’ choices across repeated instances of the same task (recall, four random tasks were repeated for each subject) which give us direct, subject-level measures of *behavioral noise*. This analysis therefore relies on new data, not reported in the previous analysis.

We estimate the model using Bayesian hierarchical techniques, which optimally combine individual-level information with group-level evidence (Gelman et al. 2014). This allows us to study distributions of individual-level parameters based on relatively few decision tasks (details and code are provided in Online Appendix E). We normalize the variance of the prior to  $\sigma = 1$  throughout, so that sampling variance is measured relative to the variance of the prior,  $\nu/\sigma$ . This is done without loss of generality and to improve comparability across studies, simply leading to a rescaling of the equation (see Natenzon 2019 for an equivalent simplification).<sup>31</sup> We execute tests on distributional differences and correlations in individual-level parameters based on the means of the individual-level posteriors throughout. All comparisons are within-subject, leveraging our two stage design, unless specified otherwise. We report four main findings:

First, we find that, conditional on the information subjects have about probabilities, estimates of  $\gamma$  indicate strong (and similar) levels of sampling variance in DfD and DfE, with  $\gamma$  estimates well below the unbiased benchmark of 1. To estimate  $\gamma$  in a way that makes DfE and DfD estimates comparable, we estimate the model in DfE on the actually experienced probabilities (i.e., probabilities implied by the sample subjects have drawn), rather than the lottery’s true probabilities.<sup>32</sup> Because of this, we must make an assumption on how subjects

---

<sup>31</sup>We estimate the model on choice data while leveraging our within-subject design. That is, we estimate the model using the data from both treatments, and assuming that the parameters governing the prior remain the same across the two treatments, while leaving the other model parameters free to vary. This allows us to maximize the informative content of our sparse choice stimuli. See Online Appendix E for details.

<sup>32</sup>We assume throughout that the initial Beta parameters, before any samples are observed, are  $\alpha = \beta = 0.1$ . This assumption derives from our general inference framework, based on a diffuse Dirichlet space – see Online Appendix A.1 for details. While values smaller than 1 are plausible (they imply that subjects expect relatively few outcomes in our general inference framework), our results are not sensitive to variations of this value within that range.

perceive the log cost-benefit ratio in cases in which the subject fails to sample both lottery outcomes before making a choice. Panel A in Figure 11 shows the cumulative distribution function of individual-level  $\gamma$  estimates under the assumption that DMs are “naive” in the sense that they judge costs and benefits to be equal in such cases. In panel B, we instead assume DMs are sophisticated in the sense that they realize that larger log-odds imply larger log cost-benefits; the correlation measuring the degree of sophistication thus must be estimated as an endogenous parameter (see Online Appendix E for details and additional results).

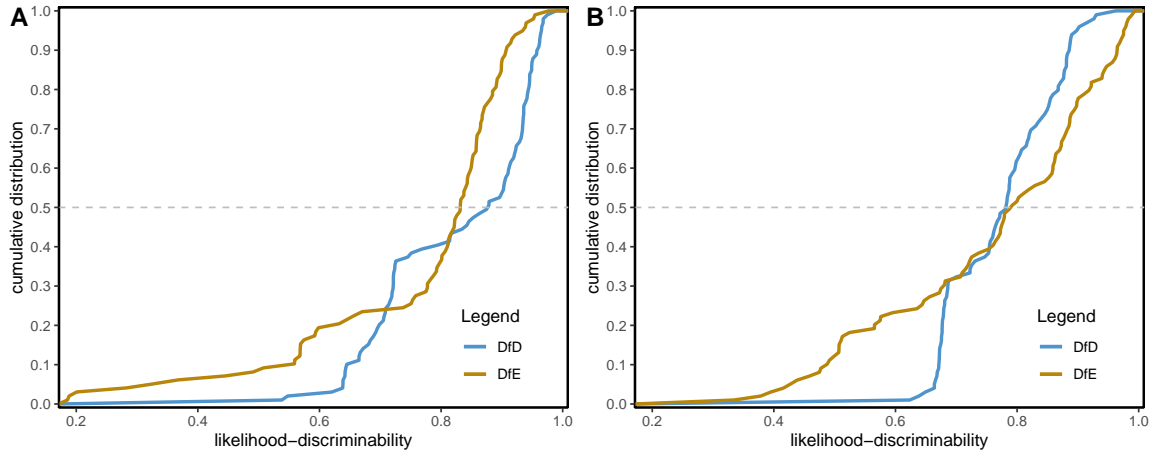


Figure 11: Structural estimates, DfD versus DfE

The figure shows structural estimates of the model parameters. Panel A compares likelihood-discriminability  $\gamma$  in DfD and DfE for a naive decision maker, who assumes costs and benefits to be equal when one of the outcomes has not been observed. Panel B compares likelihood-discriminability,  $\gamma$ , for a sophisticated DM, who (correctly) infers that log-odds and log costs-benefits are correlated in the choice problems. The correlation coefficient is thereby estimated endogenously from the data (see Online Appendix E for details).

Regardless of the approach taken, two findings stand out from Figure 11. First, in both DfD and DfE,  $\gamma$  falls well below the unbiased benchmark of 1, suggesting a strong role for inference bias in both settings as predicted by our model. Second, the distributions of  $\gamma$  estimates are similar in both DfD and DfE.<sup>33</sup> This is important because our model explains the GAP between these settings not via differences in  $\gamma$  but rather via the very different effects the model predicts  $\gamma$  has in DfD vs. DfE environments. The results therefore assure us that the model parsimoniously explains differences in lottery choices across treatments, conditional on the information available to subjects.

Second, we show that forced sampling in DfD and DfE results in a sharp increase in  $\gamma$

<sup>33</sup>For the naive estimates pictured in panel A, likelihood-discriminability  $\gamma$  is somewhat smaller in DfE than in DfD ( $p = 0.006$ ). For the sophisticated estimates in panel B, the two distributions produce roughly equal deviations above and below 0.5, and are not significantly different ( $p = 0.979$ ).

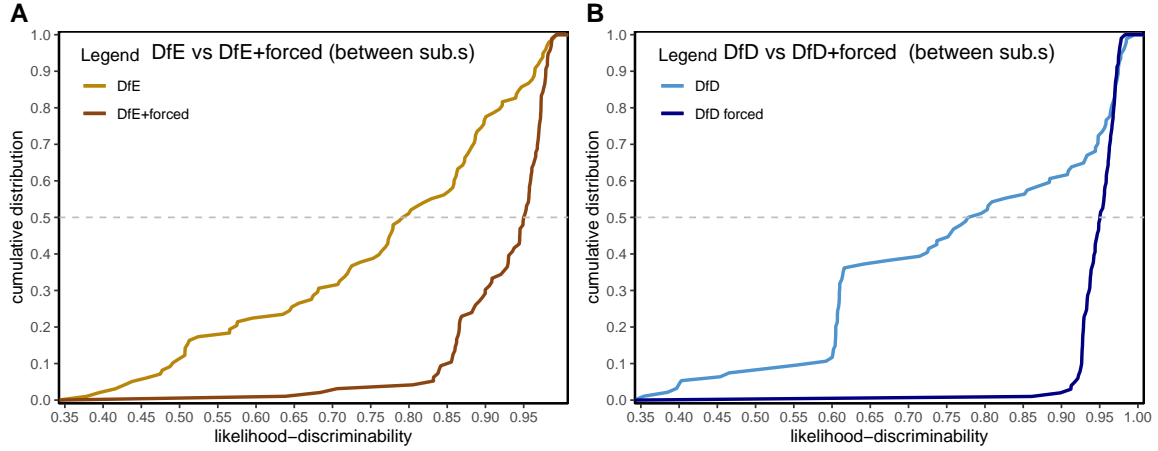


Figure 12: Structural estimates, DfD vs DfD+forced and DfE vs DfE+forced

The figure shows structural estimates of likelihood-discriminability  $\gamma$ . Panel A compares likelihood-discriminability in DfE and DfE+forced. Panel B compares likelihood-discriminability in DfD and DfD+forced.

towards 1 (the unbiased benchmark), suggesting that the intervention influences behavior (as predicted by the model) by severely reducing sampling variance and with it scope for sampling error. In panel A and B of Figure 12 respectively we plot CDFs of estimated individual-level mean  $\gamma$  estimates in DfE<sup>34</sup> and DfD with and without forced sampling.<sup>35</sup> In both cases, forced sampling causes a sharp rightward shift in the  $\gamma$  parameter, with medians in both cases of about 0.95 suggesting a near elimination of coding noise and inference bias.<sup>36</sup>

Third, we show that forced sampling in DfD and DfE – which, recall, caused a convergence in behavior between the two treatments – also causes a convergence in  $\gamma$ . This suggests (as our model predicts) a causal linkage between the two findings: joint convergence of  $\gamma$  in the two treatments towards 1 (signalling the disappearance of inference bias) causes lottery choice patterns to converge, suggesting (as predicted by the model) that coding noise was responsible for their initial divergence. Panel A of Figure 13 directly compares  $\gamma$  in DfD+forced and DfE+forced. Over most of the distribution, the panel shows that discriminability converges across the two treatments, suggesting that subjects are similarly free of inference bias in the two settings – a finding that matches the similar revealed risk

<sup>34</sup>In DfE we plot estimates that assume subjects make sophisticated inferences about the cost-benefit ratio, as discussed above.

<sup>35</sup>For this analysis, we use a between-subject comparison in both cases since DfE vs DfE+forced can only be compared between subjects; in DfD, replacing this with within-subject comparisons yields very similar results (cfr. Online Appendix E).

<sup>36</sup>Estimates also reveal a sharp reduction in cross-subject variance. This too is a prediction of the model, since the treatment is predicted to have similar impacts on both initially high and low noise subjects.

aversion in choices in the two settings. Indeed, non-parametric tests detect no significant difference between the two distributions ( $p = 0.376$ ).<sup>37</sup>

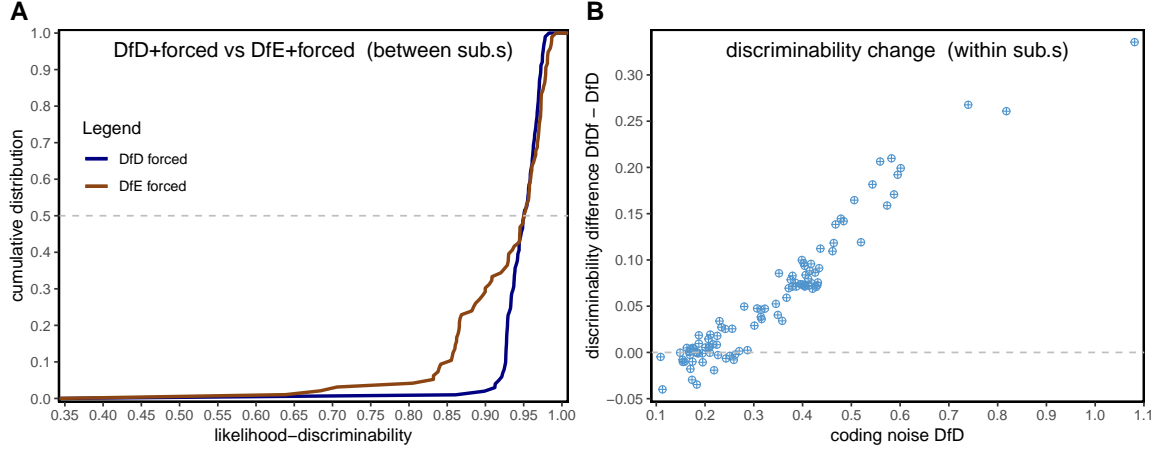


Figure 13: Effects of forced sampling, Structural estimates

The figure shows structural estimates of the model parameters. Panel A directly compares likelihood-discriminability  $\gamma$  in DfD+forced and DfE+forced. Panel B compares likelihood-discriminability,  $\gamma$ , in DfD without and with forced sampling. Panel B plots coding noise in first stage DfD against the change in likelihood-discriminability when forced sampling is introduced.

Finally, Panel B of Figure 13 illustrates the reason for this effect by plotting coding noise  $\nu$  (measured in the DfD choices in Stage 1 of the experiment) against the difference between  $\gamma$  in DfD and DfD+forced (defined as  $\gamma_2 - \gamma_1$ , with subscripts indicating the stage of the experiment), exploiting our within-subject design. The figure shows clearly that the effect of sampling is most pronounced for those subjects who had the largest coding noise to begin with. These results strongly support an additional prediction of the model: that sampling should have the strongest effect on subjects who have relatively high coding noise to start with (i.e., relatively small ‘spike counts’  $\hat{\alpha}$  and  $\hat{\beta}$ ). This is a consequence of the fact that the reduction in coding noise decreases at a decreasing rate with further samples. The figure thus shows in a particularly sharp way how strong the effect of forced sampling is on likelihood-discriminability in the DfD treatment.

<sup>37</sup>Nonetheless, as is clear from the graph, discriminability is somewhat lower in the left hand tail of the DfE distribution. We hypothesize that this is due to limitations on subjects’ memory, highlighting the value to subjects of having an explicit description of the outcomes and probabilities on the screen (in DfD+forced) to guard against inattention and working memory limitations.

## E Additional results

### Additional results on free sampling in DfE

Subjects take relatively few samples in our experiment, something that may be explained by the high opportunity costs faced by subjects on Prolific, who—contrary to students in lab or classroom settings—can leave as soon as they are done with the experiment and move on to other earning opportunities. The average number of samples taken is 8, which puts our study at roughly the first tercile of the distribution summarized in the meta-analysis of Wulff et al. (2018). Samples taken, however, generally tend to be lower in tasks comparing lotteries with sure outcomes, as we use here. The average subject on the average task takes 3.3 samples from the safe option, but 4.3 samples from the risky option. However, samples vary greatly between individuals, ranging from 2 on average (1 per option) to about 40.

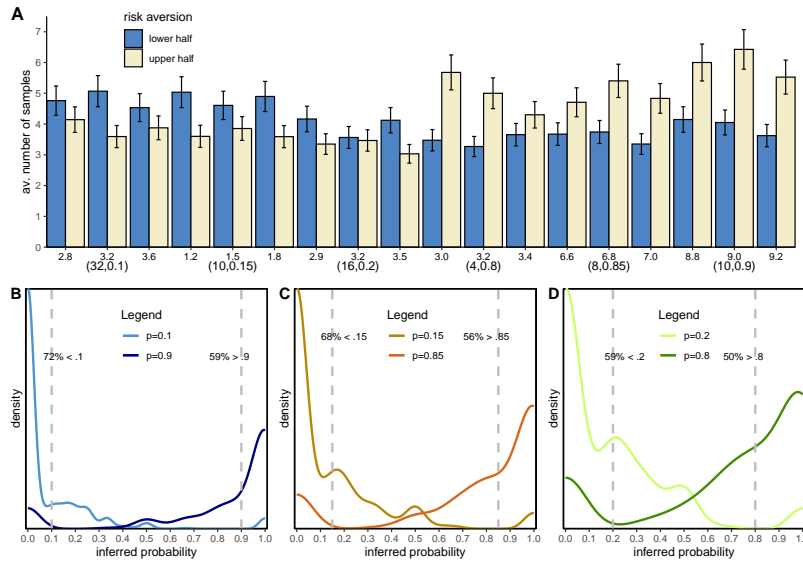


Figure 14: Samples by probability and risk aversion

The figure shows the number of samples taken from the risky option by probability and risk aversion at the task level resolution in Panel A. Risk aversion is assessed as the proportion of safe choice in the first, DfD part of the experiment, after removing repeated tasks. The categorization is obtained using a median split. Error bars show  $\pm 1$  standard error. Panels B through D show the distribution of sampled probabilities by different actual probabilities.

Panel A in Figure 5 examines the average samples by probability from the risky option at the task resolution. The samples are presented following a median split on risk aversion in the first, description-based, part of the treatment, implemented as the proportion of choices of the sure amount. This aims to test our model prediction according to which samples should vary with the underlying probability depending on the initial risk aversion of the DM. These



predictions are strongly supported by the evidence presented in the figure. Risk averse DMs take few samples from small-probability lotteries, but sample significantly more from large-probability lotteries. For the least risk averse half of the sample, we observe a (somewhat weaker) trend in the opposite direction. This aligns with our prediction, according to which risk averse DMs should have less of a conflict between noise and sampling bias in small probability lotteries, thus reaching a decision more quickly.

The small number of samples taken is reflected in the probabilities people experience. This is illustrated figure 14, panels B through D, which plot distributions of probabilities inferred from the actual samples a DM observed. For small probability lotteries, subjects experience a smaller probability than the true one in 66% of cases overall, while getting a correct picture in some 3.4% of cases. For large probability lotteries this picture is reversed, with 55% of samples over-estimating the true probability, and only 2.2% resulting in a correct estimate. The asymmetry we see between small and large probabilities suggests that the larger samples taken for large probabilities result in a more balanced picture.

### **Nonparametric within-subject results**

Here, we replicate the nonparametric between-subject analysis in the paper by presenting within-subject comparisons wherever this is possible. The descriptions of the figures are self-contained.

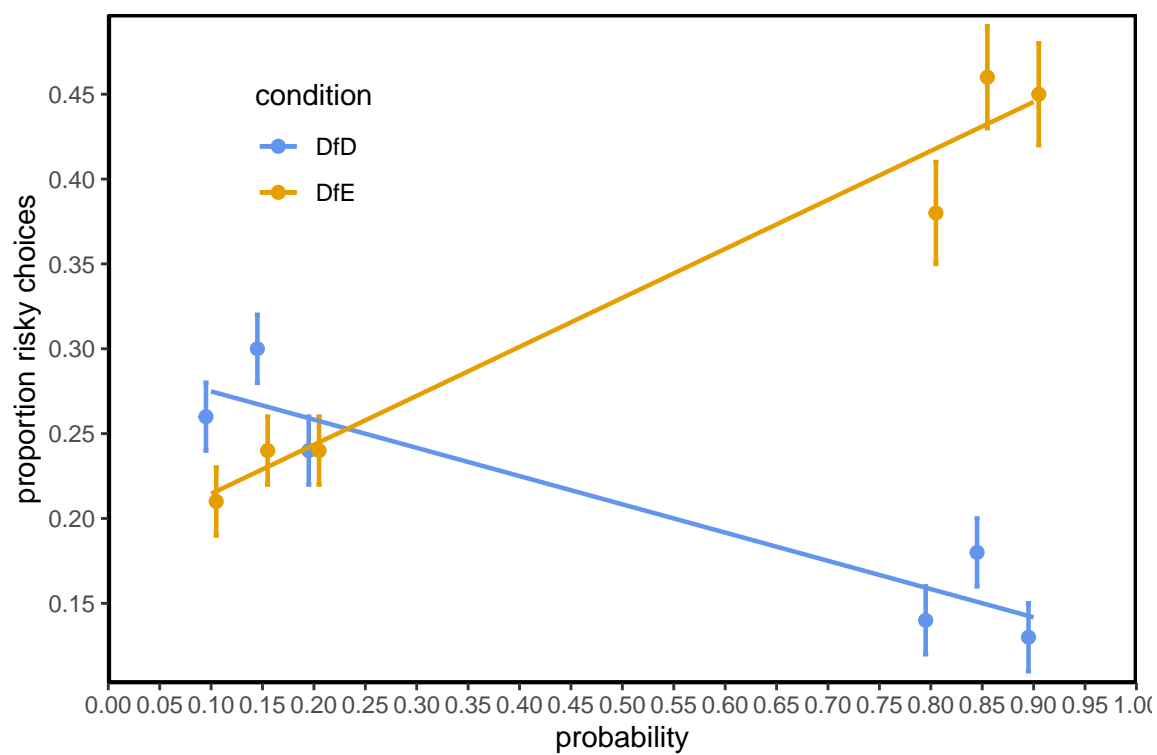


Figure 15: The GAP: within-subject

Choice proportions by probability for the decision-experience gap: DfD versus DfE. Error bars indicate 1 standard error.

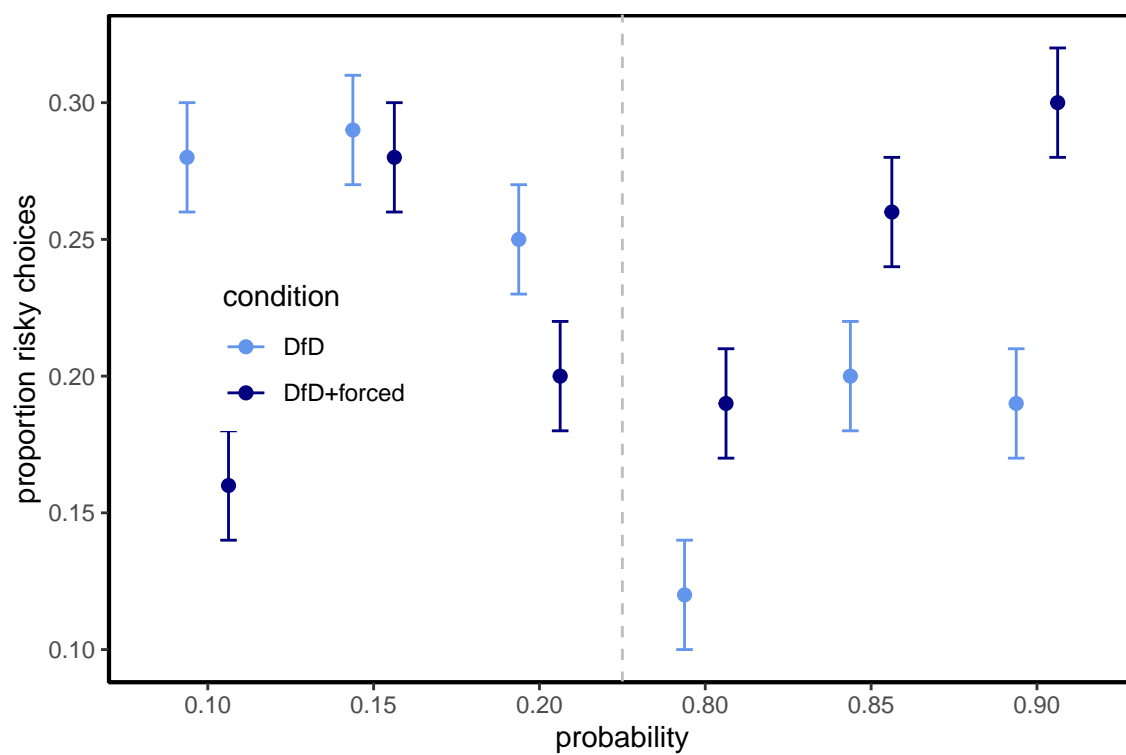


Figure 16: DfD+forced vs DfD within subject

Choice proportions by probability, within-subject comparison between DfD+forced and DfD. Error bars indicate 1 standard error.

## Figures at task level

Here, we show all figures for which we averaged across  $c$  at the probability level at a task-level resolution. The figure descriptions are self-contained.

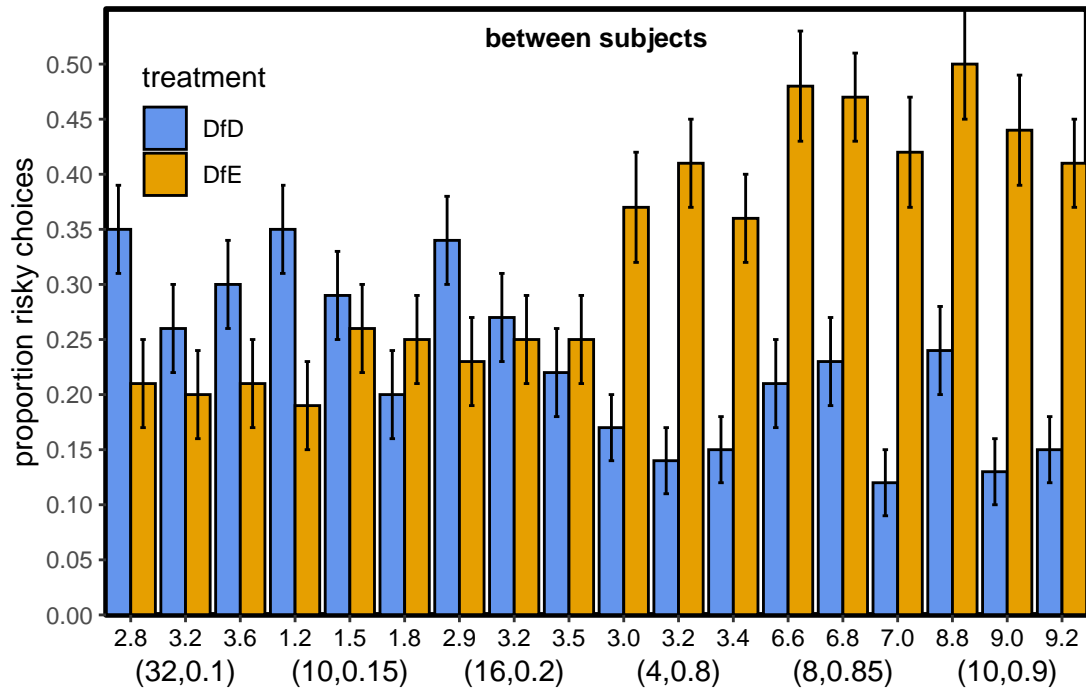


Figure 17: The GAP at the task level (between-subjects)

Choice proportions by task for the decision-experience gap: DfD versus DfE. Error bars indicate 1 standard error.

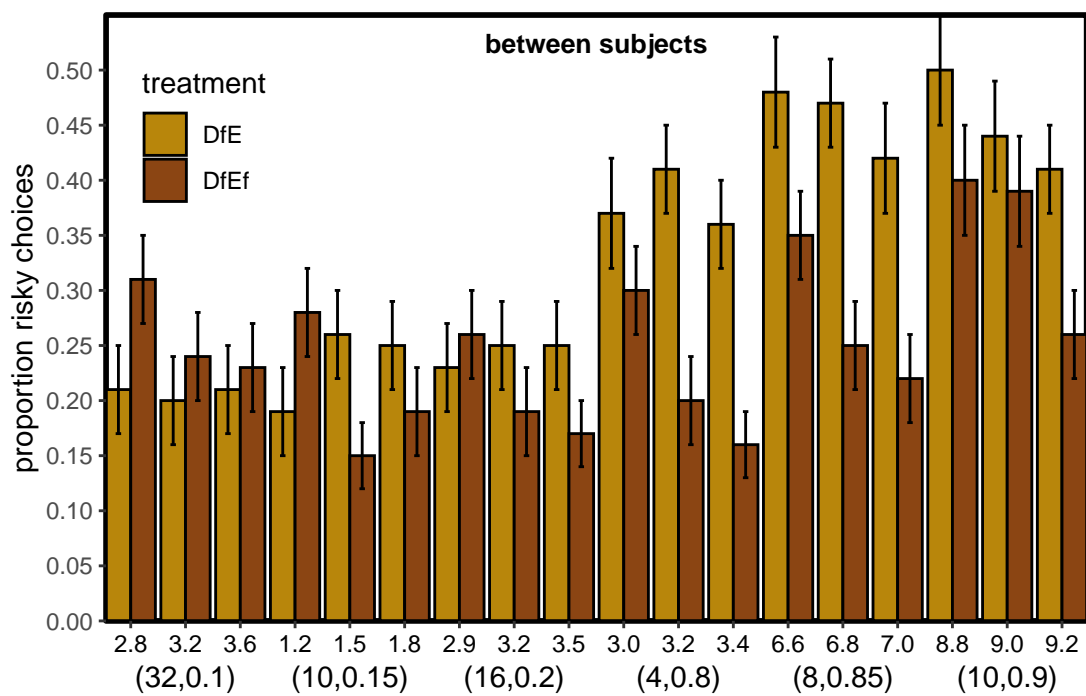


Figure 18: DfE+forced versus DfE at the task level (between-subjects)

Choice proportions by task for DfE+forced compared to DfE. This comparison is only possible between-subjects. Error bars indicate 1 standard error.

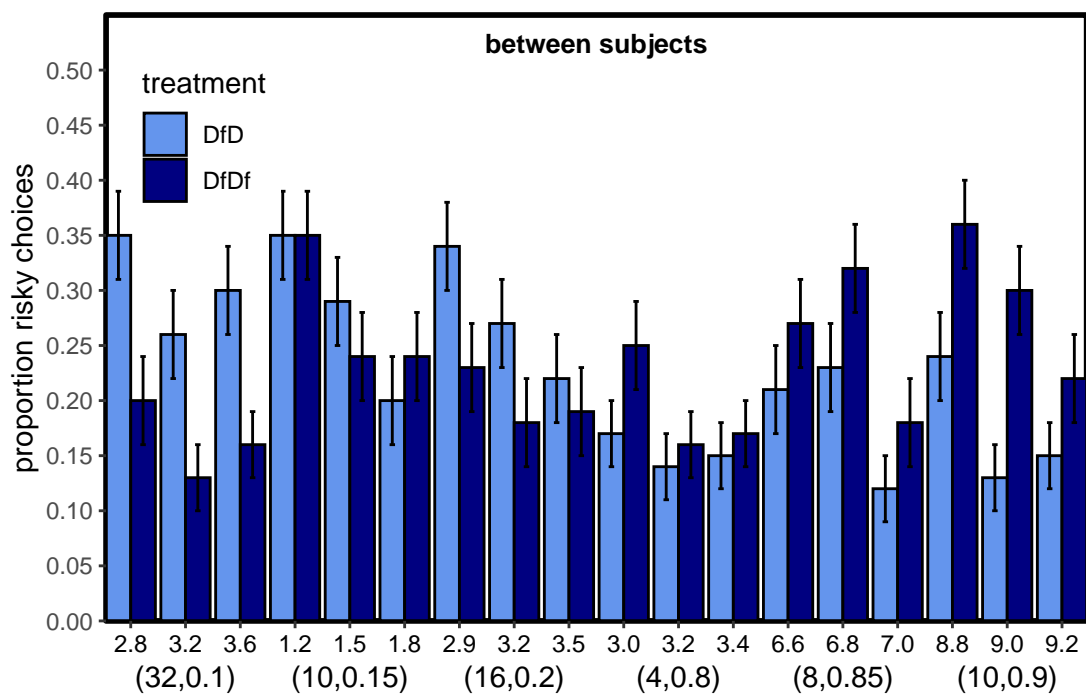


Figure 19: DfD+forced versus DfD at the task level (between-subjects)

Choice proportions by task for DfD+forced compared to DfD. Error bars indicate 1 standard error.

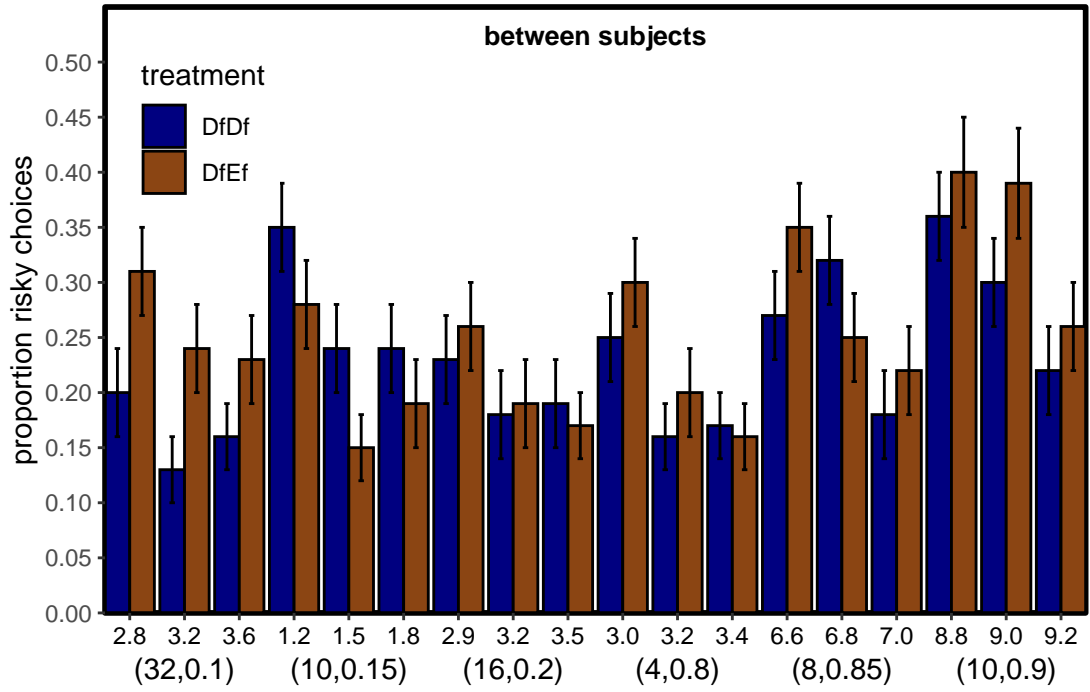


Figure 20: DfD+forced versus DfD at the task level (between-subjects)

Choice proportions by task for DfD+forced compared to DfD. Error bars indicate 1 standard error.

## Within-subject structural results

This section contains within-subject structural comparisons for those cases where we used between-subject comparisons in the main text, but within-subject comparisons are possible. The descriptions of the graphs are self-contained.

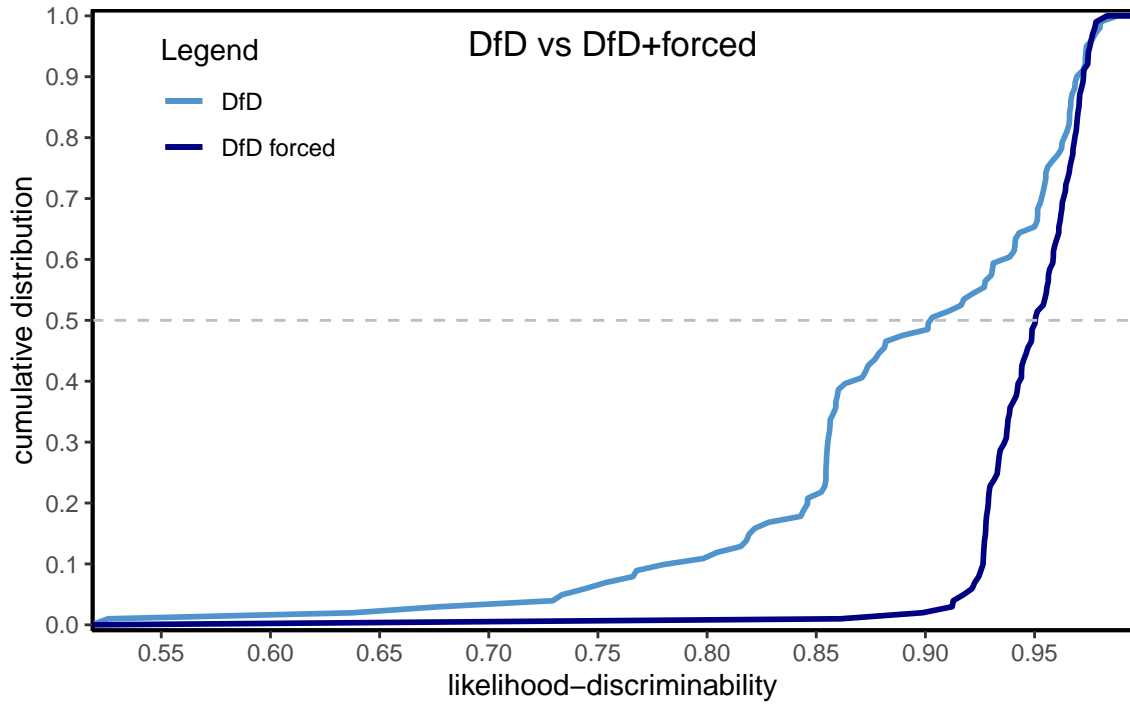


Figure 21: Likelihood-discriminability in DfD vs DfD+forced, within subject

Likelihood-discriminability,  $\gamma$ , empirical cumulative distribution function of individual-level posterior means. Within-subject comparison between DfD and DfD+forced.

## F Instructions to Subjects

### F.1 Stage 1 Instructions

Subjects in all treatments, were given the following instructions prior to Stage 1.

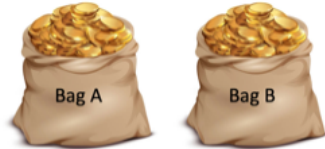
#### Instructions: Bonus

Please pay close attention to the following instructions. We will ask you **comprehension questions** about the instructions. **Anyone who answers these questions correctly the first time will receive a \$0.25 bonus.**



## Part 1 Instructions: Digital Bags

1. There will be two Parts to this experiment.



2. Part 1 will consist of **several Tasks**. In each Task you will choose between two **digital bags** -- Bag A and Bag B
3. Each bag contains **20 coins** and each coin is worth some amount of money to you as a **bonus**.

| Bag A  | Bag B                 |
|--|-----------------------|
| 80% are worth \$2.00<br>20% are worth \$0.00 | 100% are worth \$1.00 |

Example: In the example above, 80% of the coins in Bag A (i.e. 16 coins) are worth \$2, while 20% of the coins (4 coins) are worth \$0. On the other hand, 100% of the coins in Bag B are worth \$1.

4. No coin in any bag is worth more than \$35.
5. We will **randomly digitally draw one coin** from one of the two bags (Bag A or Bag B), and use that coin to determine how much money to add to your bonus. Each coin in the bag is **equally likely** to be drawn.
6. Your job is to decide **which bag** you would like us to randomly draw a coin from for your payment, by clicking one of the two buttons as in the example below.

Make Your Choice

|                       |                       |
|-----------------------|-----------------------|
| Choose Bag A          | Choose Bag B          |
| <input type="radio"/> | <input type="radio"/> |

Example: In the earlier example, if you choose Bag A there is an 80% chance you earn \$2 and a 20% chance you earn \$0. However, if you choose Bag B there is a 100% chance you earn \$1.

Please answer the following comprehension questions about the following pair of bags:

Bag A

70% are worth \$3.00  
30% are worth \$0.00

Bag B

100% are worth \$2.00

If you answer all of these questions correctly **on the first try** we will pay you a bonus of \$0.25.

In the example above, what is the likelihood (percentage chance) of earning exactly **\$3** if you choose **Bag A**.

0%

30%

70%

100%

In the example above, what is the likelihood (percentage chance) of earning exactly **\$2** if you choose **Bag A**.

0%

30%

70%

100%

In the example above, what is the likelihood (percentage chance) of earning exactly **\$3** if you choose **Bag B**.

0%

30%

70%

100%

In the example above, what is the likelihood (percentage chance) of earning exactly **\$2** if you choose **Bag B**.

0%

30%

70%

100%

### Instructions: Details

1. We will give you a total of **22 tasks** in Part 1. In each task, the contents of the bags will be **different**.
2. At the end of the experiment, we will **randomly select 10% of participants** to actually be paid a bonus based on their choices.
3. If you are selected to be paid a bonus, we will randomly select one of the tasks and use your choice to determine your bonus.

## F.2 Stage 2 Instructions

In Stage 2, subjects assigned to the DfD treatment were given the following instructions:

### Part 2 Instructions

1. The choices in Part 2 will be similar to the choices in Part 1.
2. We will give you a total of **22 tasks** in part 2. In each task, the contents of the bags will be **different**.
3. At the end of the experiment, we will **randomly select 10% of participants** to actually be paid a bonus based on their choices.
4. If you are selected to be paid a bonus, we will randomly select one of the tasks and use your choice to determine your bonus.

Subjects assigned to DfE or DfE+forced were initially given the following instructions:

### Part 2 Instructions

In Part 2 tasks, you will be making the same kind of choices you made in Part 1. However, unlike in Part 1, in Part 2 we will not describe what is contained in each bag. Instead you can learn about the contents of the bags by **sampling coins from them**.

Subjects assigned to DfD+forced or DfD+forced were initially given the following instructions:

### Part 2 Instructions

In Part 2 tasks, you will be making the same kind of choices you made in Part 1. However, you will also be allowed to **sample coins from each bag** before making your choices.

After this, subjects in DfE or DfD+free were given the following instructions:

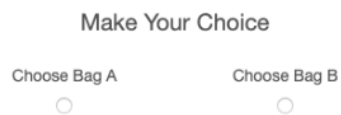
## Part 2 Instructions: Sampling

1. In this Part, in order to help you make your decision, we will allow you to **"Sample"** **from each of the bags**. We will show you buttons like the ones below. Each time you click on a button, it will **draw one of the coins** from the corresponding bag and show you how much is on it. This **won't affect your earnings** -- it is just a chance to learn about each bag.



Example: In the example above, you have clicked bag A and the computer randomly drew a coin worth \$2 from it (shown in green).

2. You can Sample from each bag **as many times as you like**. Each time you do, the computer will "put the coin back in the bag" before you sample again.
3. When you are finished sampling, just click on a button like the ones below to make your real choice (the choice that actually affects your earnings). The computer will then randomly draw one of the 20 coins from the bag to determine your bonus.



while subjects in DfE+forced or DfD+forced were instead given the following instructions:

### Part 2 Instructions: Sampling

1. In this Part, in order to help you make your decision, we will allow you to **"Sample"** **from each of the bags**. We will show you buttons like the ones below. Each time you click on a button, it will **draw one of the coins** from the corresponding bag and show you how much is on it. This **won't affect your earnings** -- it is just a chance to learn about each bag.



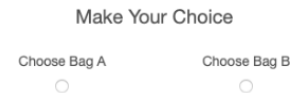
Example: In the example above, you have clicked bag A and the computer randomly drew a coin worth \$2 from it (shown in green).

2. You must Sample from **each bag 20 times**, drawing each of the 20 coins out of each bag. Each time you sample, the computer will take the sampled coin out of the bag before you sample again.



Example: In the example above, you have sampled 8 times so far from Bag A and 6 times so far from Bag B. You must sample a total of 20 times from each Bag before you can make your real decision.

3. When you are finished sampling, just click on a button like the ones below to make your real choice (the choice that actually affects your earnings). The computer will then randomly draw one of the 20 coins from the bag to determine your bonus.



Finally, all subjects were given these instructions prior to the beginning of Stage 2:

### Part 2 Instructions: Details

1. We will give you a total of **22 tasks** in part 2. In each task, the contents of the bags will be **different**.
2. At the end of the experiment, we will **randomly select 10% of participants** to actually be paid a bonus based on their choices.
3. If you are selected to be paid a bonus, we will randomly select one of the tasks and use your choice to determine your bonus.