

Belief Elicitation and Behavioral Incentive Compatibility

David Danz, Lise Vesterlund, Alistair J. Wilson*

February 2022

Abstract: Subjective beliefs are crucial for economic inference, yet behavior can challenge the elicitation. We propose that belief elicitation should be incentive compatible not only theoretically but also in a de facto behavioral sense. To demonstrate, we show that the binarized scoring rule, a state-of-the-art elicitation, violates two weak conditions for behavioral incentive compatibility: (i) within the elicitation, information on the incentives increases deviations from truthful reporting; and (ii) in a pure choice over the set of incentives, most deviate from the theorized maximizer. Moreover, we document that deviations are systematic and center-biased, and that the elicited beliefs substantially distort inference.

Keywords: Incentive compatibility, belief elicitation, binarized scoring rule, experiments.

* All authors: Department of Economics, 230 Bouquet Street, Pittsburgh, PA. Danz: University of Pittsburgh, danz@pitt.edu. Vesterlund: University of Pittsburgh and NBER, vesterlund@pitt.edu. Wilson: alistair@pitt.edu. We thank Stefano DellaVigna and four thoughtful reviewers for their constructive and very helpful feedback. We also thank Tim Cason, Yoram Halevy, Paul J. Healy, Steffen Huck, Alex Imas, Dorothea Kübler, Ryan Oprea, Isabel Trevino, Emanuel Vespa, and seminar audiences at Chicago-Berkeley, Stanford, Vibes, Vienna University of Economics and Business, UCSB, and WZB Berlin. We thank Felipe Araujo, Prottoy Akbar, Mallory Avery, Conor Brown, Ying Kai Huang, Matthew Raffensberger, Yuriy Podvysotskiy, and Tianyi Wang for help with the design and implementation of the experiment, and for excellent research assistance Pun Kanatip and Marissa Lepper.

1. INTRODUCTION

Information on individual beliefs is central to our ability to draw inference on economic decisions (Manski, 2004). Absent data on what people think and expect, we are often unable to discriminate between alternative models of choice, gauge the limits of rationality, or test new equilibrium concepts. While individual beliefs are of clear importance, eliciting them is not straightforward.

Researchers incentivize individuals to truthfully report their beliefs—with incentive compatible rules outperforming incompatible ones (Nelson and Bessler, 1989; Palfrey and Wang, 2009; Schotter and Trevino, 2014) and these in turn dominating unincentivized elicitation (Gächter and Renner, 2010; Wang, 2011; Trautmann and van de Kuilen, 2015).¹ While the empirical evidence demonstrates the advantage of incentive compatible elicitations, it is less evident that the offered incentives succeed in eliciting the participant’s true belief. The task at hand is complex. To identify beliefs through revealed choice it must be that for every possible belief there is a corresponding choice in the mechanism that uniquely maximizes the agent’s objective.

Critical for designing elicitations is that the agent’s objective is accounted for. Early elicitations were incentive compatible only for individuals who were risk-neutral expected-utility (EU) maximizers. However, an observed tendency to report beliefs toward the center was eventually seen as evidence that the assumption of risk neutrality was misguided, and that risk-averse participants were reporting more-conservative beliefs than those they truly held.² Recent design efforts have therefore aimed to find elicitations that make truth telling theoretically incentive compatible for a broader set of preferences.

We argue that to secure truthful revelation, elicitation mechanisms need to not only be incentive compatible in a purely theoretical sense, but also in a behavioral one. We propose for assessment two weak conditions for behavioral incentive compatibility, that information on deployed incentives increases truthful revelation; and that most participants, when given a choice over the pure incentives, select the outcome thought to be uniquely maximizing under the mechanism (i.e., a requirement of behavioral incentive compatibility for a representative agent).

To demonstrate we explore a state-of-the-art belief elicitation, the binarized scoring rule (Hossain and Okui, 2013, henceforth BSR). The BSR is seen as a particularly promising alternative to elicitations requiring risk neutrality because its incentive compatibility expands to arbitrary EU preferences—in fact, to any decision maker who maximizes the overall chance of winning a prize. Building on the insights of Roth and Malouf (1979), this is achieved by linking reported beliefs to a pair of state-contingent

¹ Performance measures for belief elicitation mechanisms include the distance between reported beliefs and realized outcomes or objective probabilities, assessment of the dispersion, support, and additivity of a report distribution, and internal consistency between beliefs and actions. We focus our analysis on elicited beliefs over an unambiguous objective prior, and assess performance by the distance between reports and the given objective prior.

² For recent reviews see Schotter and Trevino (2014) and Schlag, Tremewan and van der Weele (2015).

lotteries, where for each distinct belief, the mechanism provides a lottery pair with a stochastically dominant reduction. That is, decision makers who maximize their chance of winning are given strict incentives under the BSR to reveal their true belief.

In addition to being incentive compatible for a wider set of preferences, initial empirical evidence shows that the BSR outperforms its narrower forerunner, the quadratic scoring rule (Hossain and Okui, 2013; Harrison and Phillips, 2014). Weakened theoretical assumptions and evidence for superior relative performance has quickly rendered the BSR the preferred elicitation.³ However, limited evidence exists on whether subjects behave in a truth-telling manner, and the conservative reporting patterns that identified failures in quadratic-scoring elicitations have also been detected in BSR elicitations. For example, in Babcock et al. (2017), despite the qualitative comparative statics for beliefs mirroring behavior, the elicited reports appeared overly conservative.

We examine whether the incentives offered in the BSR lead to truthful revelation and in doing so demonstrate the need for elicitations to satisfy weak conditions for behavioral incentive compatibility. We start by asking whether the BSR incentives—in particular, participants’ knowledge of the precise quantitative incentives being offered—encourage truthful reporting. The main finding is that information on the offered incentives increases false reports and causes systematic bias toward the center. Later, we directly assess the BSR incentives and find that most participants, when given a choice, fail to select the outcome assumed to be uniquely maximizing under the mechanism. Finally, we demonstrate the substantial impact on inference of using the center-biased reports under the mechanism and point to the potential mistakes from using elicitations that violate weak conditions on behavioral incentive compatibility.

We start by eliciting beliefs in four distinct treatments, solely varying participants’ information over the BSR incentives. A challenge for examining whether information on the mechanism’s quantitative incentives encourages truth telling is that we do not know participants’ true beliefs. Hence, to assess the consequences of providing information on the BSR incentives we begin by evaluating reports over an objective prior. While our study also elicits posteriors that exhibit similar patterns, our focus is on eliciting the induced prior probability—where we know what a well-incentivized participant should report.

³ Recent applications of the BSR include studies on gender and coordination (Babcock et al., 2017), investment and portfolio choice (Hillenbrand and Schmelzer, 2017; Drerup et al., 2017), coordination (Masiliūnas, 2017), matching markets (Chen and He, 2017; Dagnies et al., 2019; Sonsino, Lahav, and Levkowitz, 2020), biased information processing (Graeber, 2020; Hossain and Okui, 2019; Erkal et al., 2021; Rafkin, Shreekumar, and Vautrey (2021), cheap talk (Meloso et al., 2018), risk taking (Ahrens and Bosch-Rosa, 2019), information source choice (Charness, Oprea, and Yuksel, 2021), memory and uncertainty (Enke, Schwerter, and Zimmermann, 2020; Enke and Graeber, 2021), discrimination (Aksoy, Chadd, and Koh, 2021; Dianat, Echenique, and Yariv, 2018; Koutout, 2020), correlated and motivated beliefs (Hossain and Okui, 2020; Oprea and Yuksel, 2020; Cason, Sharma, and Vadovič, 2020), auctions with private information (Corazzini, Galavotti, and Valbonesi, 2019); information and market behavior (Filippin and Mantovani, 2019; Renes and Visser, 2019); institutions and cognitive ability (Choi et al., 2020); overconfidence and bargaining (Colzani and Santos-Pinto, 2020); second-order beliefs (Dustan, Koutout, and Leo (2020); beliefs in strategic games (Aoyagi, Frechette, and Yuksel, 2021; Castillo et al., 2019).

Our first treatment, a baseline *Information* treatment, provides transparent quantitative information on the incentives. Participants are told that their chance of winning is maximized by truthful reporting, and provided with a description of how the BSR mechanism is implemented, simple numerical information on the offered lotteries for all provisional responses, and end-of-period feedback. The Information treatment reveals frequent and substantial deviations in the reports from the objective prior—where we term any deviation from the induced prior as false.⁴ Further, the large rate of false reports is systematically center-biased—with false reports more likely for non-centered than centered priors, and the direction of deviations ‘pulling-to-center.’ The remaining three treatments help distinguish whether the center-biased reporting results from confusion over the reporting task or from the BSR incentives used in the elicitation. Cognitive limitations, hedging motives, and flat incentives are all factors that could result in center-biased reporting under the BSR.⁵

To demonstrate, consider a participant in our experiment asked to guess which of two urns was randomly selected (a Red urn or a Blue urn). Given an objective prior that the Red urn is selected with probability π_0 , participants are asked to report their belief q , incentivized by a state-contingent lottery pair: a lottery with a $1 - (1 - q)^2$ chance of winning \$8 if the Red urn was selected; and a lottery with a $1 - q^2$ chance of winning the \$8 if the Blue urn was selected.

Table 1. BSR Incentive Lotteries

Submitted Belief on Red	Chance of receiving \$8 if:	
	Urn is Red	Urn is Blue
1.0	100%	0%
0.9	99%	19%
0.8	96%	36%
0.7	91%	51%
0.6	84%	64%
0.5	75%	75%

Table 1 illustrates a subset of the state-contingent lotteries offered under the BSR. Consider a participant with an induced prior of 0.8. Truthfully reporting $q = 0.8$ yields a 96 percent chance of winning the prize if the selected urn is Red, and a 36 percent chance if the urn is Blue. Reporting the induced prior therefore secures an expected chance of winning of 91 percent, the largest feasible percentage. But cognitive limitations or hedging

⁴ As such, we term reports of the objective prior as truthful. This true/false terminology is chosen for clarity, and does not imply that all participants are assumed to understand that the objective prior is the true likelihood.

⁵ We consider hedging in a broad sense, as in Dean and Ortoleva (2017).

motives can draw participants to deviate. For example, reporting a more-conservative belief of $q = 0.7$, decreases the chance of winning by 5 percentage points conditional on Red (from 96 to 91 percent) and increases it by 15 percentage points conditional on Blue (36 to 51 percent). By design, for the modeled decision-maker who perfectly reduces lotteries, the tradeoffs across the two lotteries do not warrant misreporting. While a movement towards the center from the true chance on Red of $\pi_0 = 0.8$ yields a larger improvement on the less likely Blue-urn lottery, this increase is for the modeled decision maker not large enough to compensate for the losses on the more likely Red-urn lottery. However, decision makers motivated by hedging or who struggle to reduce the compound lottery could be drawn to more-centered reports. Further this draw to center could intensify with the certainty of the belief. For example, consider a participant reporting 0.89 instead of a true belief of 0.99 on Red. Doing so increases the chance of winning conditional on Blue by an order of magnitude (from 2 to 21 percent) but decreases the chance of winning conditional on Red by a tiny amount (from 99.99 to 98.79 percent). Importantly, the cost of deviations from the theorized maximizer are trivially small for those tempted by a more-centered report.⁶

To examine whether the BSR incentives give rise to center-biased reporting in the Information treatment, we deploy two treatment modifications: one supplementing the provided quantitative incentive information and the other eliminating it. In the first modification, the *Reduction-of-compound-lottery* (RCL) treatment, we provide participants with a calculator to help them reduce the chance of winning on the lottery-pair offered. While the RCL treatment helps reduce the rate of false reports and eliminates the pull-to-center effect, false reports continue to arise at high rates and remain more likely on non-centered priors.⁷ In sharp contrast, when removing all quantitative information on the BSR incentives in our *No-Information* treatment, we find a substantially lower rate of false reports and no systematic bias in reporting. Our comparison of our Information versus No-Information treatments (henceforth Information-No-Information) shows that participants understand the task at hand—as they report the induced priors at high rates in the absence of quantitative information on incentives—and that center-biased reporting results from the provided information on incentives.

Next, we conduct a *Feedback* treatment to further explore the perverse finding that information on the BSR incentives substantially and systematically biases reports. While our Information-No-Information comparison reveals the effect of incentive-information between-subjects, the Feedback treatment instead evaluates the within-subject response using a more incremental revelation of incentive information (through end-of-round

⁶ The cost of reporting a prior other than that given is a 1 percentage point decrease in the chance of winning for a 10 percentage-point deviation in the report. For an induced prior of 0.8 a truthful report secures a 91 percent overall chance of winning. The chance of winning drops to 90 percent with a false report of 0.7 and to 88 percent with a false report of 0.6.

⁷ The RCL treatment points to divergence of the reported belief from the known prior, partially resulting from subjects' misconception of the incentives (à la Cason and Plot, 2014).

feedback on the earned lotteries). We find that false reports start out at the same low-rate as the No-Information treatment, but that the rate increases as participants gradually acquire information on the quantitative incentives, eventually reaching the level of false reports in the Information treatment. Confirming our initial finding, we see that it is information on incentives that leads to center-biased reporting.

To assess the impact of our finding we review prior applications of the BSR and find that the majority provide information on incentives with the potential of eliciting center-biased reports. In understanding the impact on inference of such reporting, we first explore a simple theoretical model, and second examine the magnitude of the inferential effect by replicating the paper by Niederle and Vesterlund (2007, henceforth NV) on gender and competition. We replicate the study using the BSR to elicit subjective beliefs over tournament ranking, using treatments with and without quantitative information on incentives. Through our Information-No-Information design we therefore evaluate the impact of information on incentives and the potential impact of center-biased reporting.

Abstracting from the elicited beliefs, both replications mirror the original NV-finding that men compete more than women. However, treatment differences are uncovered when drawing inference that factor in the beliefs. The *NV-No-Information* treatment fully replicates the original findings that men are more overconfident than women and that this gap in confidence helps explain the gender gap in tournament entry. Looking instead at the *NV-Information* treatment, the elicited beliefs appear center-biased, and (in opposition to the original results) we find both that there is *no* gender gap in confidence and that confidence *does not* help explain the gender gap in tournament entry. Thus, our replication demonstrates that information on the BSR incentives can lead to a qualitatively distinct conclusion—here in a setting that can affect policies aimed to address the advancement of women.

Our results demonstrate the substantial inferential impacts from using an elicitation that fails a weak condition for behavioral incentive compatibility: that information on the incentives increases truthful reporting. A similarly weak condition for behavioral incentive compatibility is that most participants, when given a direct choice over the set of incentives under the mechanism, will select the assumed maximizing outcome. While our Information-No-Information comparison tests the first condition (compatibility within the elicitation environment), the second condition can be assessed in an *Incentives-only* treatment, in which subjects face the same incentives as in the BSR treatments—absent the elicitation framing. Participants in our Incentives-only treatment are shown the set of available lottery pairs under the BSR, mirroring the options in Table 1, and asked to pick their preferred pair under a fixed probability for the Red event. As additional evidence that the elicitation is not behaviorally incentive compatible, we find that the large majority of participants fail to select the outcome assumed to be the unique maximizer. Deviations instead move towards the center, consistent with the center-bias documented in the elicitation, and indicative of incentives as the root cause of center-biased reporting.

The paper is structured as follows: Section 2 presents our design and results from the Information treatment where participants get detailed information on the quantitative incentives. Section 3 contrasts the Information results with those of two treatments, one that removes all information about the incentives, and another that adds a calculator to help reduce the compound lottery inherent to the BSR. Section 4 focuses on gradual release of information on the quantitative incentives through end-of-period feedback. Section 5 shows that our main results extend to the elicited posterior beliefs, and Section 6 examines the impact of biased elicitation on inference and the impact on belief elicitation more generally.

2. BASELINE DESIGN AND RESULTS

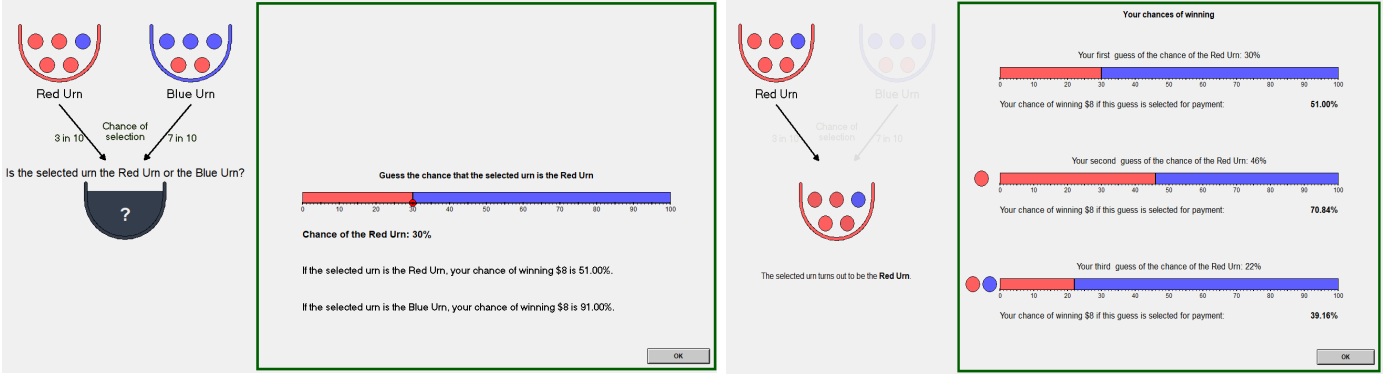
We start by evaluating whether information on the incentives offered under the BSR elicitation increases truthful reporting. We conduct four initial treatments to assess elicited beliefs of an objective prior and to explore the response to information on incentives. The environment is held constant across all four treatments. Undergraduate students were recruited to participate in an individual decision-making task at the Pittsburgh Experimental Economics Laboratory (PEEL).⁸ Each treatment consisted of three separate sessions (with a recruitment aim of 20 participants per session). The procedures of the experiment, the number of periods, the elicited belief scenarios, as well as the offered incentives were all held constant across treatments. In terms of exposition, we describe the common features of the experimental environment as we introduce the Information treatment.

2.1. Information Treatment

The *Information treatment* is designed as a baseline: an implementation of the BSR that presents participants with clear information on the quantitative incentives associated with any provisionally considered report on the belief. An Information treatment session (like all the treatments we examine) consists of ten periods, where each period has three sequential decisions. Participants are paid for one of the three decisions in two of the ten periods.⁹

⁸ In addition, subjects had to be 18+ years old to be eligible for participation. Invitations to all sessions were gender balanced.

⁹ The experimental interface is programmed in z-Tree (Fischbacher, 2007). Participants received printed instructions that were read out loud and summarized in a short, scripted presentation at the start of each session (see Appendix B for instructions for the Information Treatment, with exact language-deltas for all other treatments; cf. reporting best practices articulated in De Quidt et al. 2019; and Appendix C for the slides and script used in the summary presentation). Across all treatments the average duration of a session was 71 minutes with average earnings of \$20.08, including an \$8 show-up fee.



(A) Choice Interface

(B) End-of-period feedback

Figure 1. Interface Screenshots

Within the session participants make choices in the interfaces shown in Figure 1. Panel A shows the decision screen and panel B the end-of-period feedback. At the beginning of each period participants are shown two urns, one Red and one Blue. Each urn contains five colored balls (either red or blue) where the Red Urn contains more red balls than the Blue Urn. One of the two urns is selected at random and the main task for participants is to guess the likelihood that the selected urn is Red. Participants are informed of the composition of both urns and of the prior probability π_0 that the Red Urn is selected, presented as an X-in-10 chance (see panel A). Given this information, participants are asked to submit three sequential guesses on the chance that the selected urn is Red. Guess 1 is made without any additional signals, and Guesses 2 and 3 are made, respectively, after observing the colors of a first and then a second independent draw from the selected urn.¹⁰

The decision screen in Panel A of Figure 1 shows a marked 30 percent guess on the Red Urn, secured by placing a cursor on a slider ranging between 0 and 100 percent (with one percentage-point increments). Each provisionally marked guess leads to an offered state-contingent lottery pair displayed on the screen: one if the selected urn is the Red Urn, another if the selected urn is the Blue Urn. Both lotteries are over a prize of \$8 if won, and \$0 otherwise. As noted in the introduction, the BSR incentive given a stated probability of q on the Red Urn offers a $1 - (1 - q)^2$ chance of winning the \$8 if the Red Urn is drawn, and $1 - q^2$ if the Blue Urn is drawn. Thus, the chance of winning is maximized by reporting the given likelihood Red is selected, which for Guess 1 corresponds to the induced prior, π_0 , and for Guesses 2 and 3 the Bayesian posteriors, updated in response to the draws from the selected urn.¹¹

¹⁰ Our task resembles the standard Bayesian updating task with the addition of the prior elicitation, see Benjamin (2019) for a recent review of the literature on belief updating.

¹¹ The evaluated scenarios and random realizations are held constant across treatments. Within each session all 20 participants see the same sequence of 10 scenarios, though in different random orders. While the signal realizations vary across participants and sessions of a treatment, the sequencing of scenarios considered (state, signal realizations, and sequence) is held matched across treatments.

Later we will introduce (sequentially) three additional treatments that vary the information on the quantitative incentives provided to participants. However, we first outline our results in the Information treatment, which provides clear incentive information through four channels: (i) The instructions explicitly provide the qualitative information that truthful reporting is a dominant strategy (a common feature to the presentation in all of our treatments) stating that “[t]he payment rule is designed so that you can secure the largest chance of winning the prize by reporting your most-accurate guess.” This statement is also emphasized in a slide presentation summarizing the instructions where it is the last thing participants see before making their first decision (see final slide of Appendix C). (ii) The written instructions provide a concise verbal description of how the state-contingent lotteries determine the prize realization.¹² (iii) Within the interface, as participants move their provisional belief, the screen is instantly updated to provide clear quantitative information on the state-contingent probabilities of winning. This can be seen in Panel A of Figure 1 in the two lines below the input slider. With $q = 0.3$ selected, the interface displays the associated chances of winning the prize for each realization of the selected urn; in this case 51 percent if Red, and 91 percent if Blue. (iv) Finally, as shown in Panel B of Figure 1, participants receive feedback information on the selected urn at the end of each period, as well as the realized quantitative chance of winning the \$8 prize given the state realization (the selected urn) and their submitted Guesses 1 through 3.

After the ten periods (30 elicitations total) we elicit risk attitudes (encoded as switch points on pricelists) and ask participants to respond to a Cognitive Reflection Test (Frederick, 2005).¹³ One participant per session is randomly selected to be paid for these end-of-experiment elicitations. Finally, participants complete a post-experiment questionnaire on demographics, and provide a self-assessment of their comprehension of the incentives and the extent to which they reported their most-accurate guess.

2.2. Information Treatment Results

In examining whether BSR incentives secure truthful reporting we focus our analyses on Guess 1, the elicitation of the induced prior. The induced prior is unambiguous and should be reported back by every participant who understands the offered incentives and seeks to maximize their chance of winning, independent of the participants ability to

¹² The explanation of how the chance of winning was determined in the state-contingent lotteries relied on a comparison of the reported guess to that of two (uniform) random numbers, thereby avoiding the presentation of formulas, or the understanding of a squared error (see Wilson and Vespa, 2018).

¹³ In each row of the first table, participants chose between a sure payoff of \$4 and a lottery $p \cdot \$8 \oplus (1 - p) \cdot \0 with p increasing in each row from 0 to 1 in steps of 0.1 (see Bruner, 2009). In the second table the lottery was the same in all rows with $\frac{1}{2} \cdot \$8 \oplus \frac{1}{2} \cdot \0 and the sure payoff increased over the rows from \$0 to \$8 in steps of \$0.80 (see Abdellaoui et al., 2011). Tables 3 and 4 were the same as 1 and 2, respectively, except that all prize payoffs were scaled by a factor of 1.25.

Bayesian update.¹⁴ Section 4 presents a complementary analysis for Guesses 2 and 3, showing that the same qualitative results hold.

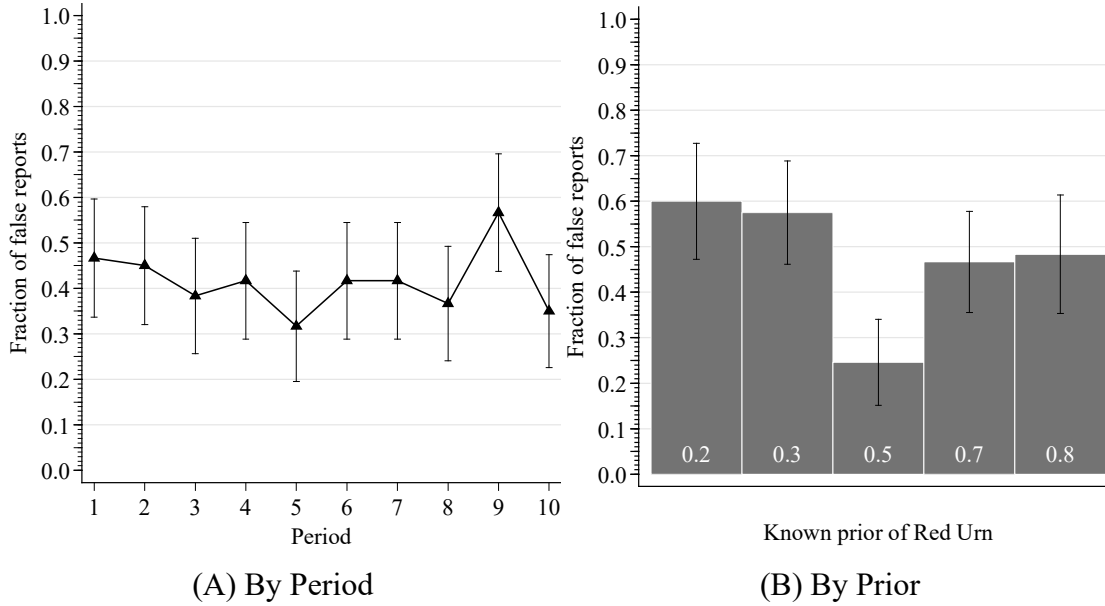


Figure 2. False-report rate in Information Treatment.

Figure 2 illustrates the rate of false reports (any elicited belief q that differs from the induced prior π_0) by period across the session, and by the objective prior. Panel A shows a substantial rate of false reporting over the ten periods, averaging 41.5 percent, which is maintained without a time trend across the session ($p=0.842$).¹⁵ False reports are widespread, with 85 percent of participants failing to report the induced prior in one or more of the ten periods. These deviations from truth telling are particularly concerning considering the incentive compatibility for a general family of underlying preferences, the qualitative statement that truth telling will maximize participants' chances of winning, and the prior evidence on the comparative superiority of the BSR. Panel B illustrates how the rate of false reports varies with the location of the induced prior. For non-centered priors ($\pi_0 \in \{0.2, 0.3, 0.7, 0.8\}$) we find that false reports are the norm (52.8 percent), while they are significantly less likely to occur for the exact-centered prior ($\pi_0 = 0.5$, with a 24.6 percent false report rate, $p < 0.001$ from an OLS regression with participant-clustered standard errors).

¹⁴ While Offerman et al (2009) use the elicitation of induced priors for *ex post* corrections, we use it to assess the elicitation procedure itself, see also Hao and Houser (2012) and Holt and Smith (2016). The simplistic elicitation eliminates belief formation and may help participants focus on the incentives provided (see for example, Avoyan and Schotter, 2020 on shared attention).

¹⁵ Tests of time trends are based on probit regressions of false reports on period with participant-clustered standard errors.

Conditional on a false report, the average deviation from the prior is 0.167, and the large proportion of false reports are not driven by small mistakes.¹⁶ Further, false reports tend to pull-to-center. Among the false reports for non-centered priors we find that 53.7 percent lie between the objective prior and the exact center (a stated report of $q=\frac{1}{2}$), while only 32.6 percent fall between the objective prior and the nearest extreme (with the remaining 13.7 percent of misreports being somewhere between the exact center and the distant extreme). This reveals a significant pull-to-center, with greater false reports toward the center than the nearest extreme ($p=0.058$).¹⁷ In contrast to belief-elicitation mechanisms such as the QSR, where the ‘pull-to-center’ effect was interpreted as resulting from risk aversion, the effect here is unexpected as the BSR is incentive compatible for arbitrary risk preferences. Indeed, we find no evidence that risk aversion is the culprit: Individual risk attitudes do not predict the propensity to deviate from the true prior, nor the conservative reporting tendency.¹⁸

Why are so many participants misreporting the prior? Is it the task at hand that participants fail to understand or do they simply object to reporting the given prior? Alternatively, are the offered incentives causing participants to distort their reports? We argued in the introduction that along with the miniscule deviation cost, center-biased reporting may result because participants fail to reduce compound lotteries, or because participants hold hedging motives (preferring a substantial increase in winning on the unlikely event in return for a slight decrease in the chance of winning on the likely event).

In examining the drivers of the large degree of false-reporting in our Information treatment we consider three potential channels: (i) false reports inherent to the task of reporting the prior, for example, resulting from confusion about what is being asked; (ii) false reports driven by a failure to reduce the compound lottery inherent to the BSR incentive; (iii) false reports driven by other features of the BSR incentives, for example, through a non-EU hedging motive.¹⁹ While the first channel captures confusion inherent

¹⁶The results are the same when shifting our binary definition of a false report to allow for small errors that deviate by no more than 5 percentage points from the prior (Appendix Table A.1). For example, this less strict definition of false reports give rise to a false-reports rate of 40.6 percent for non-centered priors and 17.1 percent for centered priors (different with $p<0.001$). See Cason and Plot (2014) for a similar approach when assessing the BDM mechanism.

¹⁷ Partitioning the non-centered false reports in this manner secures that deviations toward the center and the near extreme have the same width, allowing us to fairly assess the extent to which participants deviate toward the center (as opposed to the near extreme).

¹⁸ Individual false-report rates and the extent to which these move toward the center are not significantly correlated with an individual being risk averse or loving (identified by whether willingness to pay for a 50 percent chance of winning \$8 is below or above the certainty equivalent of \$4).

¹⁹ Our after-experiment survey asked participants how they made their decisions, and the responses provide anecdotal evidence that they purposefully distort reports to secure a higher chance of winning on the less likely event, and were aware of the incentives for doing so (all responses verbatim): “*I generally erred on the side of caution when picking the urns. For example, if $x=5$, I would select 50% for the red urn. If sat $x=8$ then I would pick the red a little more opportunistically.*”; “*I kept my initial answers at 50% because you get a 75% chance of getting the \$8 anyways. Then I adjusted as I saw the different outcomes.*”; “*at first, I guessed based on probability probability the urn was picked based on the dice roll and then considered the*

to our belief elicitation, the latter two relate directly to the BSR incentives and thus point to deviations that result from failed behavioral incentive compatibility.

2.3. RCL & No-Information Design

We use two additional treatments to provide insights on the source of false reporting. Both manipulate the participants' information on the incentives. In the first, we provide *additional* information specifically tailored to limit misunderstanding of the compound lottery, providing a calculator that reduces it to a simple lottery (the *RCL treatment*). In the second, we remove all information on the quantitative incentives (the *No-Information treatment*).

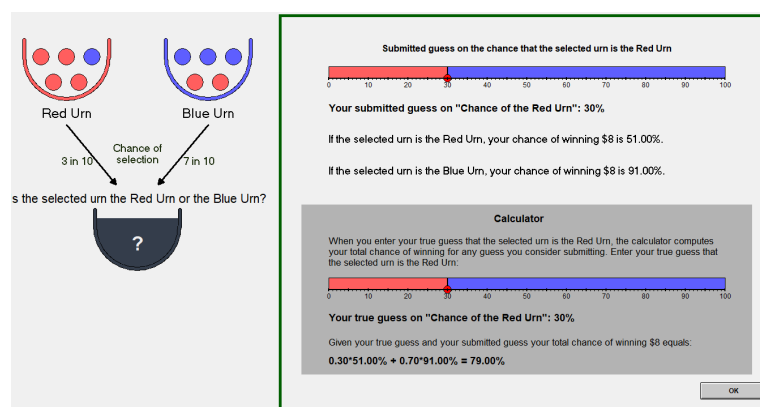


Figure 3. RCL treatment screenshot

Our Reduction of Compound Lotteries (RCL) calculator provides participants with a tool to compute the total chance of winning the \$8 prize for any report (see Figure 3 for the interface). The RCL calculator asks participants to enter their true belief and determines the overall chance of winning for each potential report. The calculator therefore helps participants verify that truth telling maximizes their chance of winning by reducing the offered lottery pair for any true and stated belief.²⁰ Beyond the addition of the RCL calculator (and supplementary instructions on how to use it) the treatment is otherwise identical to the Information treatment.

balls that were drawn from the bag; however, I quickly realized that since I am pretty risk adverse, sticking to a 50-50 chance would result in being paid the \$8 75% of the time regardless of which urn was selected. I mostly stuck to that model as I proceeded through the experiment. When i felt daring, I would move my guesses a little bit around the 50-50 mark (but never very far).”; “I believe that leaving each chance at 75% was my best chance of making the most money in the experiment.”

²⁰ Like the more-involved explanation of the BDM in Healy (2017, 2018) the hope was that this would help participants understand that the mechanism was incentive compatible. While more substantial explanations and training may have enhanced participants' comprehension of the mechanism's incentive compatibility (see, for example, Burfurd and Wilkening, 2018), because belief elicitation is typically a secondary measure in experiments, we opted for an aid that would not substantially increase the length of the instructions.

In contrast, our No-Information design holds constant all qualitative details from the Information treatment but removes all quantitative information on the incentives. Participants are still told that the procedure was designed so that truthful reporting will maximize their chances of winning and that \$8 is at stake as a prize, but they are uninformed on the chances of winning the prize.²¹ In addition to removing the description of the mechanism in the instructions, the No-Information interface also removes the numerical information on the state-contingent lotteries at each provisionally selected belief (the two lines below the input bar Panel A of Figure 1) and the end-of-period feedback on the earned chance of winning for each Guess (the three ex post probabilities in Panel B of Figure 1).²²

The RCL treatment offers a channel to assess the extent to which an inability to reduce compound lotteries is driving the Information results. Given the lack of *any* incentive information in No-Information, the level of false reports in this treatment serves to identify factors other than the incentives (for example, confusion with the task). By removing all quantitative incentive information, any difference in false reports with the Information treatment is identified as coming from some feature of the BSR incentives. The relative differences between the three treatments can therefore help to decompose the effects on false reporting.

2.4. RCL & No-Information Results

Three sessions were run for each treatment, with 60 participants in No Information and 59 in RCL.²³ Paralleling our data presentation for the Information treatment, Figure 4 reports the false-report rate by session period (panel A) and by the objective prior (panel B).²⁴ We note that while the RCL treatment reduces the frequency of false reports, the reduction is even greater in the No-Information treatment.²⁵ We also see that the pattern of greater false reporting for non-centered than centered priors is reduced but not eliminated in the RCL treatment, while it disappears entirely in No-Information.²⁶ Thus, although an improved ability to reduce compound lotteries decreases false reports, the best results are obtained by eliminating quantitative information on the BSR incentives.

²¹ Participants are informed that “[t]he precise payment rule details are available by request at the end of the experiment.” Of the 60 subjects in the treatment, only one requested this information.

²² The end-of-period feedback screen in No-Information instead provides feedback only on the realized selected urn.

²³ One RCL session under-recruited and ran with 19 participants.

²⁴ While there are significant differences in the average absolute error across treatments, the difference is driven by the rate of false reports, hence our focus on this measure. Conditional on a false report there are no significant differences in the magnitude of the deviation from truth telling, with treatment-average deviations for false reports ranging between 0.156 and 0.183.

²⁵ For comparison we note that Holt and Smith (2016) use a similar two-urn-guessing paradigm to elicit beliefs under the QSR, the BDM, and the lottery choice method. Evaluating only one prior of 0.5 they find false report rates on the prior elicitation ranging between 20 to 33 percent.

²⁶ Further the patterns in the results continue to hold when eliminating small mistakes within 5 percent of the induced prior (see Online Appendix Table A.1).

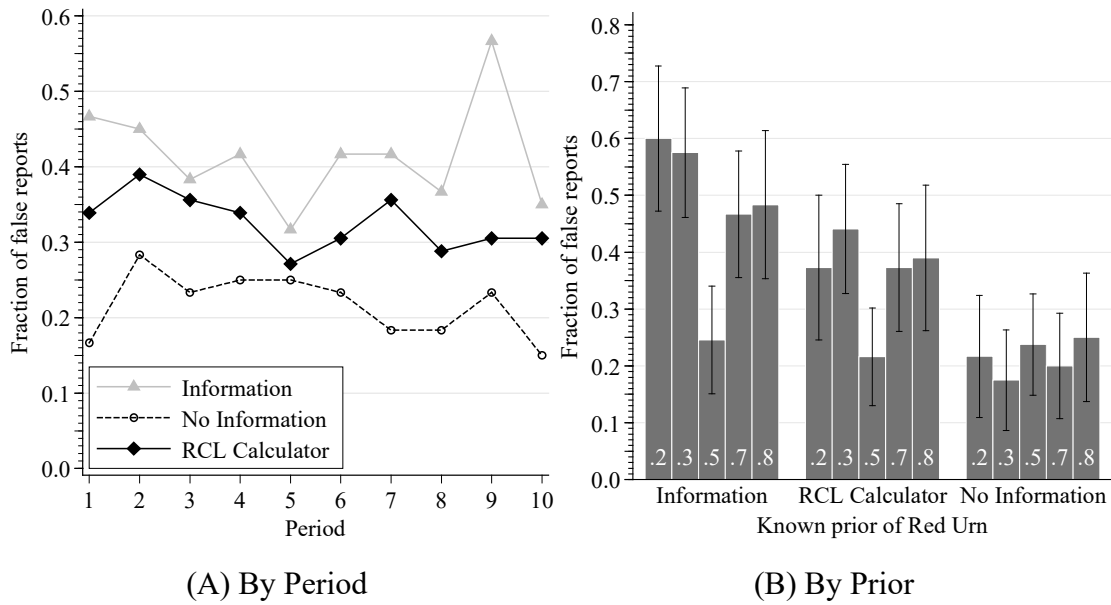


Figure 4. False-report rate in Information, No-Information and RCL treatments

As evidence of center-biased reporting we illustrate in Figure 5 the distributions of reports at asymmetric priors of 0.2 and 0.3 (panels A and B, respectively). On the left of each panel, we indicate the overall rates of truthful reporting for the asymmetric priors by treatment. On the right of each panel, we illustrate the distributions of false reports. Compared to the RCL and No-Information treatments we see for the Information treatment not only the low rate of truthful reporting, but also the center bias with greater false reporting toward the center than the near extreme.

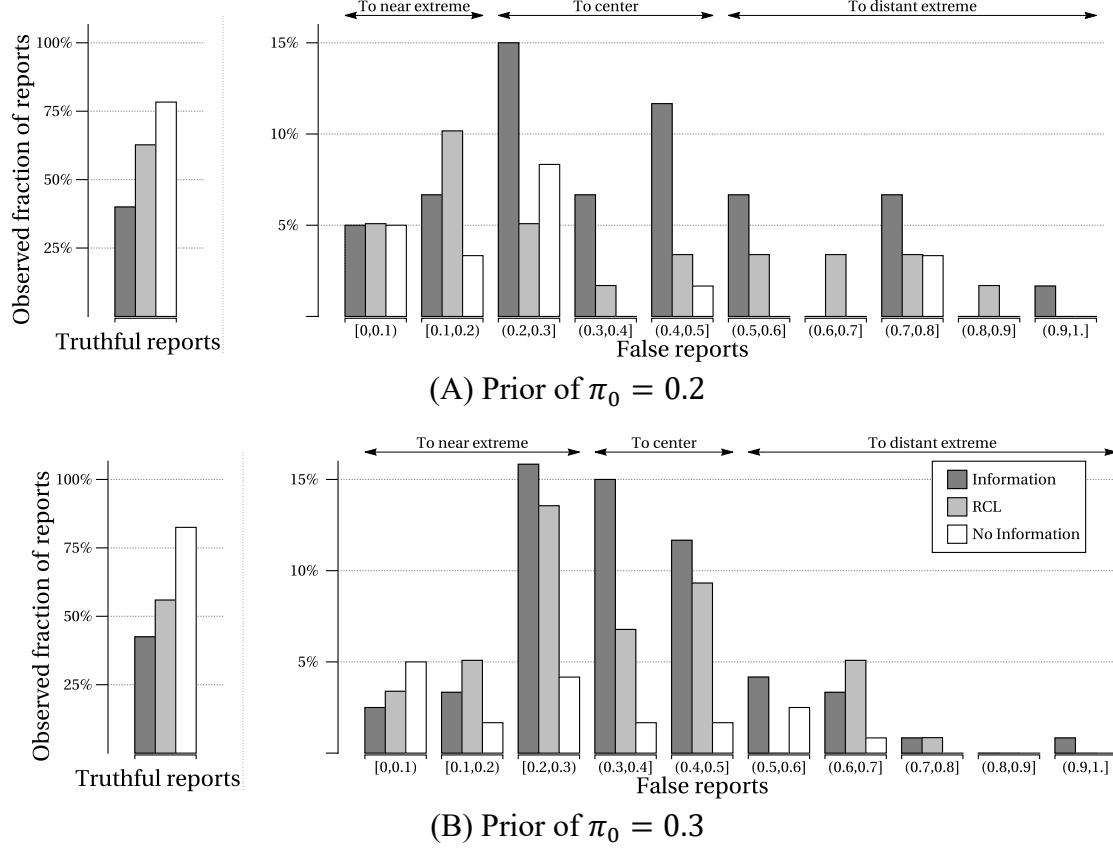


Figure 5. Proportion truthful, and distribution of response for false reports

The center-bias effects illustrated in the figures are quantified in Table 2 which shows the false-report rates with participant-clustered standard errors, where the first three rows confirm that the treatment effects illustrated in Figure 4 and 5 are significant.²⁷ The first data column in Table 2 indicates treatment-level false-report rates across all prior elicitations (pooling centered and non-centered priors).²⁸ The two *By-Prior* columns separate the false-report assessment into two subcategories: those occurring when the induced prior is centered ($\pi_0 = 0.5$), and those when it is non-centered ($\pi_0 \neq 0.5$). Finally, the last three columns decompose the false reports for non-centered priors into three regions, mirroring those labeled at the top of panels A and B in Figure 5, to assess the extent to which reported beliefs are skewed toward the center. We examine the proportion of non-centered priors for which a false report: (i) moves toward the center (false reports of $q \in (\pi_0, \frac{1}{2}]$ when $\pi_0 < \frac{1}{2}$, and of $q \in [\frac{1}{2}, \pi_0)$ when $\pi_0 > \frac{1}{2}$, respectively); (ii) moves to the nearest extreme (false reports of $q \in [0, \pi_0)$ and $q \in (\pi_0, 1]$, respectively); and (iii)

²⁷ Unless otherwise stated, all treatment comparison p -values are obtained from two-sided tests (OLS with standard errors corrected for clusters at the individual level). Probit estimates indicate almost identical quantitative marginal effects and qualitative inference, so we focus here on the easier to interpret measures.

²⁸ Table 2 also reports the average treatment levels from two additional treatments (details below).

moves between the exact center and the distant extreme (false reports of $q \in (\frac{1}{2}, 1]$ and $q \in [0, \frac{1}{2})$, respectively). The results for the baseline Information treatment (the first data row) mirror Figure 2: more than 40 percent of the reports do not equal the objective prior; where the false-report rate is significantly greater for non-centered than centered priors ($p < 0.001$); and that false reports on non-centered priors are more likely to pull-to-center than the nearest extreme ($p = 0.058$).

Table 2. False Reports by Treatment

Treatment	False Reports			False-Report Type		
	All Priors	By Prior		$(\pi_0 \neq 0.5)$		
		$\pi_0 = 0.5$	$\pi_0 \neq 0.5$	Center	Near Extreme	Distant Extreme
Information	0.415 (0.042)	0.246 (0.047)	0.528 (0.048)	0.283 (0.042)	0.172 (0.032)	0.072 (0.014)
RCL	0.325 (0.041)	0.216 (0.043)	0.398 (0.048)	0.169 (0.033)	0.164 (0.031)	0.065 (0.014)
No Information	0.217 (0.039)	0.238 (0.040)	0.203 (0.041)	0.058 (0.018)	0.108 (0.027)	0.036 (0.014)
Feedback $_{(t=1,2)}$	0.217 (0.045)	0.236 (0.060)	0.200 (0.055)	0.031 (0.021)	0.154 (0.047)	0.015 (0.015)
Feedback $_{(t=9,10)}$	0.341 (0.053)	0.255 (0.064)	0.405 (0.071)	0.087 (0.039)	0.275 (0.062)	0.043 (0.024)
Description	0.245 (0.040)	0.196 (0.039)	0.278 (0.046)	0.108 (0.027)	0.131 (0.029)	0.038 (0.012)
N	2,630	2,630		1,568		

Note: Standard errors in parentheses clustered by participant (299 clusters) recovered from three separate joint estimates on the false report proportion in the prior elicitation: (i) *All priors*, dependent variable an indicator for $q \neq \pi_0$, with treatment level estimation; (ii) *By Prior* column pair, same dependent variable as All priors, but with separate treatment estimates for centered/non-centered prior location; and (iii) *False-Report type* column triple, treatment-level estimation for the division of non-centered false reports into three mutually exclusive regions: center (between the π_0 and $\frac{1}{2}$), near extreme (between the closer extreme 0/1 and π_0), and distant extreme (the further of 0/1 and $\frac{1}{2}$).

Relative to Information we note that the RCL calculator reduces the rate of false reports. Across all elicited priors the RCL treatment leads to a 9.0 percentage point reduction in false reports. While this reduction is not significantly different from zero ($p = 0.130$), it is significant when we focus solely on the non-centered priors (a 12.9 percentage point reduction, $p = 0.056$).²⁹ Adding to the reduction in false reports is that we

²⁹ While the working hypothesis motivating the RCL treatment was one-sided—that reducing the lottery will help participants understand the incentive compatibility—we report two-sided tests for consistency.

no longer see false reports that pull-to-center.³⁰ The greater frequency of false reports on non-centered than centered priors however remains (an 18.2 percentage point difference, $p < 0.001$).

The proportion of false reports decreases even further when participants have no quantitative information on the BSR incentives (significantly lower than Information and RCL, $p < 0.001$ and $p = 0.056$, respectively).³¹ Further, there is no evidence of center-biased reporting. False-report rates are no greater on non-centered than centered priors ($p = 0.317$) and there is no pull-to-center ($p = 0.175$ for a two-sided test, where the difference has the opposite sign).³² With No Information substantially reducing the rate of false reports and eliminating center-biased reporting, we infer that both effects are causally linked to knowledge of the quantitative BSR incentives.³³ The No-Information treatment therefore demonstrates that false reporting in the Information treatment does not simply arise from confusion over the task. Participants report the objective prior at high rates, independent of its location, provided they are *uninformed* of the quantitative incentives from doing so.³⁴

Comparing the three treatments we get a sense of what drives participants to falsely report non-centered priors in the Information treatment. The false report rate on non-centered priors is 52.8 percent in Information treatment and 20.3 percent in the No-Information treatment, which suggests that 38 percent ($= 0.203/0.528$) of the false reports can be attributed to the task itself (e.g., confusion). The remaining 62 percent are directly linked to the BSR incentives, whether through an inability to reduce compound lotteries or another feature of the offered incentives. The estimated 39.8 percent false-report rate in the

³⁰ The RCL treatment's 39.8 percent false-report rate for non-centered elicitations is more evenly distributed between those that move toward the center, and those that move to the nearest extreme ($p = 0.903$). Specifically, while deviations made toward the nearest extreme in RCL continue to occur at a similar rate to Information (16.4 vs 17.2 percent, $p = 0.851$), we find a significant reduction in the false reports moving toward the center (16.9 vs. 28.3 percent, $p = 0.034$).

³¹ The rate of false reports across the No-Information sessions decreases by 4 percentage points between the last and first two periods of the treatment, though the effect is insignificant ($p = 0.420$).

³² Intriguingly, the reduction in false reports is only seen for reports that pull-to-center (28.3 in Information vs. 5.8 percent in No-Information, $p < 0.001$) and not in those toward the nearest extreme (17.2 vs. 10.8 percent, $p = 0.160$).

³³ The difference in the rate of false reports across centered/non-centered priors in RCL is highly significant ($p < 0.001$), where the frequency of pull-to-center reports is significantly larger than in No-Information ($p = 0.004$).

³⁴ The data in our post-experimental questionnaire further bolsters the case that it is the incentives that drives the false reports. Participants are asked to rate their agreement with "I always reported my most-accurate guess on the Red urn being the selected urn." Responses were collected on a 5-point Likert scale. Looking at the fraction of answers in the Strongly Agree/Agree categories, we find, 70 percent of respondents claiming they always reported their most-accurate guess in the Information treatment, and 85 percent in No-Information ($p = 0.049$, χ^2 -test of independence; see Figure A1 in the Online Appendix for further details). Self-assessment of truthful reports is (insignificantly) higher in RCL than Information (81 vs. 70 percent, $p = 0.149$). Further, while there are no differences in comprehension of mechanism between RCL and Information, participants in No-Information are less likely to report that they understood how their pay would be calculated (72 percent) and how the submitted belief affected their pay (70 percent) than participants in the Information and RCL treatments (80 and 86 percent on pay, and 83 and 86 percent on beliefs, respectively).

RCL treatment roughly suggests that failure to reduce compound lottery accounts for approximately 25 percent ($= [0.528 - 0.398] / 0.528$) while the remaining 37 percent results from some other aspect of the incentives ($= [0.398 - 0.203] / 0.528$).

We find in our three treatments little evidence that false reporting is correlated with individual characteristics. Individual risk attitudes are not predictive of the rate of false reports or of center-biased reporting in any treatment, nor are individual cognitive reflection scores predictive of false or center-biased reporting in the Information and No-Information Treatments. Only in the RCL treatment do we see evidence that cognitive reflection impacts behavior, with higher levels of cognitive reflection being predictive of fewer false reports and less center-biased reporting. This suggests that the RCL calculator is particularly helpful to those with high levels of cognitive reflection, and that even with a tool tailored to simplify decision making under transparent incentives, cognitive skills and effort are required for subjects to benefit from such aids.

3. FEEDBACK TREATMENT

Our Information-No-Information comparison presents between-subject evidence that information on the quantitative incentives drives center-biased reporting under the BSR, which is indicative of the elicitation lacking behavioral incentive compatibility. To further explore and identify the effect as coming from information on the quantitative incentives we conduct a *Feedback treatment* with 60 additional participants. Incentive information in this treatment is gradually revealed through end-of-period feedback screen (see Panel B of Figure 1). That is, in the Feedback treatment we replicate the No-Information instructions and main decision screen, but after each period's elicitation we provide participants with the Information treatment's end-of-period feedback, which informs them on the earned probability of winning as a simple lottery (given the realized state). The quantitative incentive information provided in the Feedback treatment is therefore acquired gradually as the session proceeds and is limited to the reported beliefs and realized state.

Panel A of Figure 6 indicates the false-report rate by period across the Feedback sessions. While false reports start out at the same rate as No-Information, over time the fraction of false reports increases, eventually reaching a level that is indistinguishable from that of the Information treatment. Referring to Table 2 for comparisons and inference, we find a false-report rate of 21.7 percent for the first two periods (the $\text{Feedback}_{(t=1,2)}$ row) which grows significantly ($p=0.003$) to 34.2 percent for the final two periods (the $\text{Feedback}_{(t=9,10)}$ row). Thus, feedback on the quantitative incentives increases the frequency of false reports over the session, where the starting and ending points provide a strong match to the No-Information and Information treatments, respectively. The false-report rate in the first two periods of Feedback is statistically inseparable from the overall No-Information rate ($p=1.000$) but significantly different from Information ($p=0.001$). Conversely, the false-report rate in Feedback's final two periods is significantly different

from the No-Information treatment ($p=0.060$) but inseparable from Information ($p=0.282$).³⁵

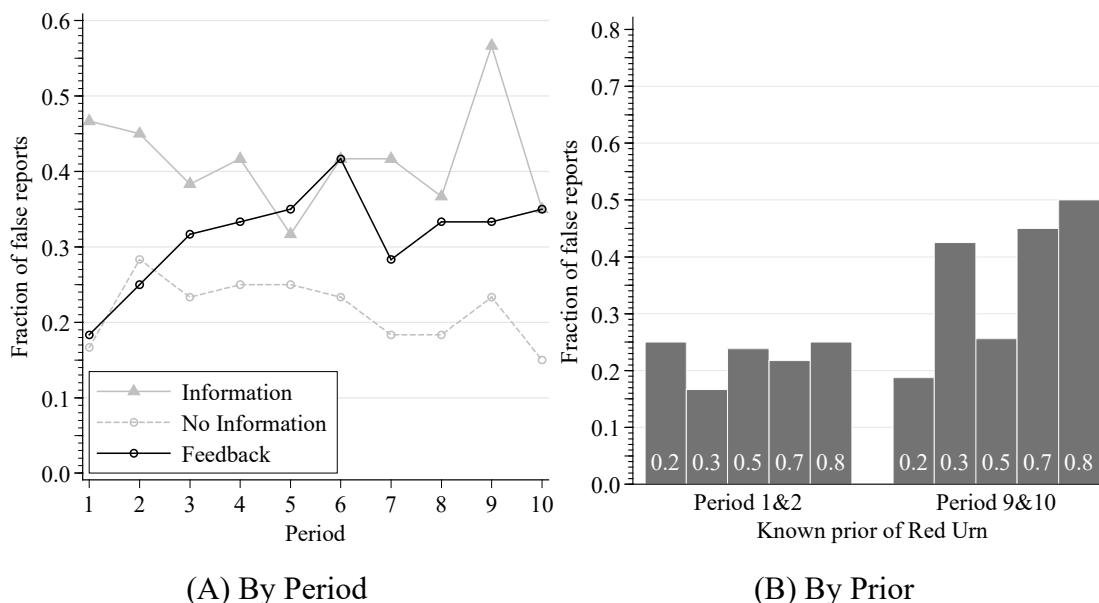


Figure 6. False-report rate in Feedback treatment

Though only provided with three *ex post* measures per period on the quantitative incentive (one for each submitted belief) the fraction of false reports in Feedback reaches the Information-treatment level within four periods. Despite the lower power when focusing on the first and last two periods, Panel B of Figure 6 suggests that participants also begin to respond differentially to non-centered and centered priors. Comparing false reports in the last two Feedback periods we find a rate of 25.5 percent for the exact-center priors, and of 40.6 percent for the non-centered priors ($p=0.089$ from a two-sided test of difference).³⁶

The Feedback treatment shows that center-biased reporting increases significantly, as participants gradually learn about the BSR incentives. These results are consistent with our Information-No-Information finding that information on incentives decreases truthful reporting—and serves as further evidence that the BSR is not incentive compatible in a behavioral sense.

³⁵ Notably, there is no significant time trend at the 5 percent level in any treatment except for the Feedback. Responses to the exit survey indicate that participants learn the incentives over time (for example, “I kept my initial answers at 50% because you get a 75% chance of getting the \$8 anyways. Then I adjusted as I saw the different outcomes”).

³⁶ Using the last three periods of data instead of the last two, the difference between non-centered and centered priors (with participant-clustered errors) is significantly different with $p=0.046$.

4. POSTERIOR REPORTS

Our analysis has focused on elicitations of the induced prior—as this provides the cleanest measure on truthful reporting, in that there is no concern over participants’ ability to perform Bayesian updating. However, it may be that quantitative-incentive information encourages accurate reporting when we are eliciting beliefs that are not objectively known, or that require effort to form.³⁷ Mirroring our analysis on the priors we examine the frequency and patterns of distortion in the reports for Guesses 2 and 3, where participants receive signals on the state and are asked to report a posterior belief.

After observing the setup for a given scenario—the composition of the two urns and the prior probability—participants report their belief on the objective prior in Guess 1. They are then sequentially shown two independent draws from the selected urn, and are asked to report an updated posterior belief after each draw. While the objective Bayesian posterior is easily determined by the analyst from the provided details, such inference is non-trivial and requires probabilistic sophistication.³⁸ As such, we expect the elicited posterior beliefs to deviate from the Bayesian posteriors. Indeed, the number of cases in which participants exactly report non-boundary Bayesian posteriors is just 6.9 percent across all treatments. Focusing on ‘truthful’ reports would therefore only capture a tiny fraction of participant decisions, and it would exclude reports of truly held, but non-Bayesian, posterior beliefs.³⁹

To assess false reporting on elicited posterior beliefs, we instead characterize reported beliefs by whether they are *distant* from the objective Bayesian posterior. We then assess the pattern of distant reports across treatments. Although the distant-posterior treatment effects are smaller than those for false-reports on the priors, the qualitative patterns mirror the previous results—in terms of the total rates, the sensitivity to location, and the pull-to-center effect. We classify distant reports as those that differ from the Bayesian posterior by more than 15 percentage points (the approximate average deviation for a false report in our prior elicitations).⁴⁰ Such reports comprise 33 percent of the data in the Information treatment, where distant reports make up 25-28 percent for the other treatments (except for the last rounds of the Feedback treatment).

³⁷ Note that a perhaps more intuitive counterargument could be made: when additional cognitive efforts is needed incentives could become less salient or less compensating for the additional effort required. Further, even if participants are taking the incentives more seriously for harder tasks, it is unclear why this should decrease rather than increase the hedged pull-to-center reporting.

³⁸ See Benjamin (2019) for a survey of the literature on behavior in Bayesian updating tasks.

³⁹ For boundary posteriors, the realized signal perfectly reveals the state through a simple inference without the need for calculation. For boundary cases where the Bayesian posterior is either 0 or 1, 84.5 percent of the elicited posteriors report the true boundary belief.

⁴⁰ In the Online Appendix, Figure A.2 shows that the results are not sensitive to what we define as ‘distant.’

Further, there is evidence of center-biased reporting when participants are informed of the quantitative incentives.⁴¹ We can further explore this by evaluating how the likelihood of reporting perfectly centered beliefs varies by treatment. In the Information treatment 11.6 percent of the elicitations with an intermediate Bayesian posterior have participants report a perfectly centered belief, where the rate in No Information is a quarter this size at 2.9 percent ($p < 0.001$ for the comparison). To illustrate the differential response between Information and No Information, Figure 7 indicates the fraction of exact-center reports in each as a function of the Bayesian posterior π (where we additionally plot the exact-center report rates for the prior elicitation). The figure makes clear that for both the prior and posterior beliefs, exact-center false reports are less likely under No Information than Information.

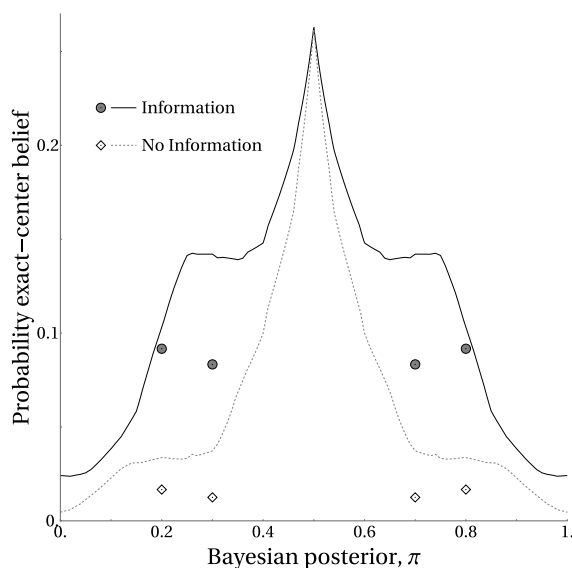


Figure 7. Exact-center reports by Bayesian posterior

In the Online Appendix (Figure A.3) we further show that information on incentives moves the range of reports toward the center. Looking across our 60 elicitation scenarios (the priors, the urn compositions, and the realized signal draws) we find a significantly greater

⁴¹ For evidence of center-biased reporting we parallel in Online Appendix Table A.2 the analysis presented in Table 2, First showing that treatment differences over distant reports are driven by elicitation where the true Bayesian posterior is in an intermediate range ($0.15 \leq \pi \leq 0.35$ or $0.65 \leq \pi \leq 0.85$): with a 41 percent distant-report rate for such posteriors in Information, decreasing to 31 percent in RCL (different from Information with $p=0.024$) and 32 percent in No Information ($p=0.040$). By contrast there are no treatment differences for posteriors in the central region ($0.35 < \pi < 0.65$). Second, showing pull-to-center by examining three distinct regions when the true posterior is *intermediate*: (i) posterior beliefs at the exact center ($q = 1/2$); (ii) posterior beliefs as the nearest extreme belief of $q = 0$ or $q = 1$ to the Bayesian posterior; and (iii) posteriors beliefs in the wrong half of the elicitation interval (so any $q > 1/2$ if $\pi < 1/2$ and vice versa). While we find no significant differences in the rates of near extreme or wrong-half posterior reports between Information and No-Information/RCL, there are strong significant differences in the frequency of exact-center reports ($p < 0.014$).

tendency for the distribution of reports to move towards the center in the Information than the No-Information Treatment.⁴²

Gaining clear identification of when a particular posterior report is distorted is certainly more challenging than for the priors, as we do not observe the participant's true updating rule. While this focuses our analysis on more-specific measures of false reporting, the general patterns mirror our findings when eliciting the induced priors. Information on the quantitative BSR incentives distorts the posterior reports toward the center.⁴³

5. DISCUSSION: IMPACT AND IMPLICATION

Our results show a large rate of false reports when providing participants with clear quantitative information on the BSR incentives. Only 15 percent of participants consistently report the induced prior under Information, in comparison to 50 percent under No-Information. Moreover, the distortions are systematically center-biased, with false reports being more likely for non-centered than centered priors and pulling-to-center. By contrast, false reports are not center biased in the absence of information on incentives.

Our consistent finding that information on incentives increases false reports points to the BSR failing a weak condition for behavioral incentive compatibility. We use this section to first discuss the potential inferential impact of such a failure, before moving on to a broader examination of the implications for belief elicitation.

5.1. Impact of center-biased reporting

Clarity of measurement focused our assessment of truth telling in the BSR on elicited reports over an objective induced prior. However, existing BSR studies mostly focus on eliciting subjective beliefs, attempting to measure a private perception where the analyst has no guidance on what the true belief might be. Should we expect our findings over the objective priors to translate to this growing number of studies using the BSR, and if so, what might the inferential impact be for those studies?

To explore these questions, we first assess whether our implementation of the elicitation mirrors those used in existing studies. Are scholars, as in our study, eliciting the likelihood of events and providing information on the quantitative incentives? Table 3 reports by study the information participants were provided, and whether an event

⁴² Using the upper (lower) quartiles for the elicited beliefs in the 58 matched scenarios where the Bayesian posterior was below (above) 0.5, we form a Wilcoxon signed-rank test for equivalence between Information and No Information. Consistent with greater movement to the center the test shows that both quartiles are significantly closer to the center in the Information treatment than in the No Information treatment ($p=0.007$ and $p=0.031$, respectively).

⁴³ Our data shows that information increases false (distant) reports both for prior and for (the more challenging) posterior reports. Further, as seen in Figure 7, the change in the exact-center reports between the Information-No-Information treatments are very similar for the prior and posterior reports. Thus, information on the BSR results in center-biased reporting—whether participants are asked to report on the (easy) prior or on the (harder) posterior.

likelihood was elicited. We find that most studies provide detailed information on the incentives, with many using two or more quantitative examples and thereby demonstrating the tradeoffs associated with center-biased reporting.

Table 3. Papers using the BSR

Paper	Full instructions	Likelihood	Dominant	Description	Quantitative Inf.	2+ Quant.Examp.	Inference
Published:							
Hossain and Okui (2013)		✓		✓	✓	✓	<i>LHS</i>
Babcock et al. (2017)	✓	✓	✓	✓	✓	✓	<i>LHS</i>
Hillenbrand and Schmelzer (2017)	✓	✓	✓	✓	✓		Both
Drerup et al. (2017)		✓	✓	✓	✓	✓	<i>RHS</i>
Masiliūnas (2017)	✓	✓			✓	✓	<i>LHS</i>
Castillo et al. (2019)	✓	✓	✓	✓	✓	✓	
Corazzini et al. (2019)	✓	✓	✓	✓	✓	✓	<i>LHS</i>
Dargnies et al. (2019)	✓		✓	✓	✓		<i>RHS</i>
Erkal et al. (2020)	✓	✓		✓	✓	✓	<i>LHS</i>
Sonsino et al. (2020)	✓	✓	✓	✓	✓	✓	<i>RHS</i>
Charness et al. (2021)	✓	✓	✓	✓			<i>LHS</i>
Rafkin et al. (2021)			✓	–	✓	–	<i>LHS</i>
Chen and He (forthcoming)	✓		✓	✓	✓		<i>RHS</i>
Oprea and Yuksel (forthcoming)	✓	✓	✓	✓			<i>LHS</i>
Working Papers:							
Hossain and Okui (2019)	✓		✓	✓			Both
Ahrens and Bosch-Rosa (2019)	✓						<i>LHS</i>
Dianat et al. (2019)		✓	✓	✓	✓	✓	<i>LHS</i>
Filippin and Mantovani (2019)		✓					<i>RHS</i>
Renes and Visser (2019)	✓	✓	✓	✓	✓	✓	<i>LHS</i>
Choi et al. (2020)	✓	✓		✓	✓	✓	<i>LHS</i>
Colzani and Santos-Pinto (2020)	✓		✓	✓	✓	✓	<i>LHS</i>
Dustan et al. (2020)	✓		✓	✓	✓	✓	<i>LHS</i>
Enke et al. (2020)	✓		✓	✓			Both
Koutout (2020)	✓		✓	✓	✓	✓	<i>RHS</i>
Meloso et al. (2020)	✓	✓	✓	✓	✓	✓	<i>LHS</i>
Aksoy et al. (2021)	✓	✓	✓	✓			<i>RHS</i>
Enke and Graeber (2021)	✓	✓	✓	✓			<i>LHS</i>
Erkal et al. (2021)	✓	✓	✓	✓	✓	✓	<i>LHS</i>
Aoyagi et al. (2021)	✓	✓	✓	✓			<i>LHS</i>
Graeber (2021)	✓	✓	✓	✓	✓	✓	<i>LHS</i>
Zimpelmann (2021)		✓	–	–	–	–	<i>RHS</i>
Totals:	81%	71%	79%	90%	70%	62%	

Note: Full Instructions – instructions available; Likelihood – elicit likelihood of an event occurring; Dominant – participants given information to reveal that truthful revelation is a dominant strategy; Description – payoffs described; Quantitative info – participants provided with some quantitative information on the incentives; 2+ Quant Example – participants provided with two or more quantitative information on incentives; Inference – whether for inference the elicited belief is used as a left or right hand side variable.

While the types of elicitations and information provided in existing BSR-founded studies suggest the potential for center-biased reporting, the impact on the inferences drawn will depend on how the elicited beliefs are used. In particular, we distinguish in Table 3 between elicited beliefs used as: (i) a left-hand-side dependent variable; and (ii) a right-hand-side independent control variable. When used as a left-hand-side variable, the effects of center bias are predictable, attenuating any estimated treatment response. In contrast, when used as a right-hand-side variable, the inferential distortions depend on the precise specification and on the relationship between beliefs and the other variables, with the bias potentially magnifying or reducing the estimated effects (see Online Appendix B.1).

To demonstrate the magnitude of the potential inferential distortions, we replicate an experimental study where elicited beliefs play a core role for inference. Specifically, we

examine the Niederle and Vesterlund (2007) study of gender and competition. We mirror the study's examination of an individual's decision to perform under a non-competitive piece-rate or a competitive tournament. However, we replace the belief elicitation at the end of the study with one of two treatments: a BSR elicitation with information on the quantitative incentives, and a BSR elicitation without that information—mirroring our Information-No-Information design.

The original NV-study finds that men more than women prefer to perform under a competitive tournament compensation, rather than under the non-competitive piece rate, and that gender differences in confidence helps explain this gender gap in competition. Our two treatments, which we label *NV-information* and *NV-no-information*, therefore serve as a testbed to evaluate first whether there is evidence of center-bias when eliciting subjective beliefs in the NV-information treatment, and second the potential inferential impact of using the elicited beliefs for analysis. Importantly, we can do this both when the beliefs are used as a dependent measure (evaluating the gender gap in confidence) and as a control (evaluating the gender gap in competition when controlling for confidence).

Design details, and elicited beliefs

Participants in our online replication of the NV-study were given two minutes to correctly add as many sets of two two-digit numbers as possible under three different performance incentives: (Task 1) a \$0.50 piece rate; (Task 2) a four-person tournament with \$2.00 per-problem-solved for the winner; and (Task 3) under the participant's preferred payment scheme (piece rate or tournament).⁴⁴ Beliefs on relative performance were then elicited at the end of the study using the BSR, with all participants given the qualitative dominance statement that the chance of winning a \$4-prize is maximized by accurately reporting the likelihood of being ranked 1st, 2nd, 3rd and 4th in their group of four. Further, participants in the NV-information treatment were given the precise conditional likelihoods of winning the \$4-prize associated with any prospective report, corresponding to our Information treatment (see instructions in Online Appendix C.3).

With the exception of the belief elicitation, our two NV treatments mirror the NV (2007) experimental design. In the original NV study participants were instead rewarded a \$1 prize for correctly guessing their rank in their group of four (1st, 2nd, 3rd, or 4th). By eliciting the modal rank the NV (2007) elicitation may be seen as less informative on the distribution of beliefs, however the coarse elicitation holds advantages that nonetheless may improve reporting over a BSR information treatment. The simple elicitation is more

⁴⁴ Standard in-person lab procedures were modified to an online format that closely mirrored those of the lab (the experimenter and other participants were visible, instructions read out loud, clarifying questions addressed in real time, see Danz et al. 2021 for the complete protocol). Participants who entered the Task-3 tournament were identified as winners if their performance in Task 3 exceeded that of the other group members in the Task-2 tournament. Participants also faced a Task-4 decision where they decided whether to submit their Task-1 performance to a piece rate or a tournament.

natural, and compared to the BSR the incentives are easier to explain (see also Abeler et al., 2019) and do not provide hedging opportunities (see Table 1). Further, inference on incentive compatibility does not hinge on the participant's ability to reduce a compound lottery. Thus, we anticipated that the coarser NV (2007) elicitation would better reflect the participants' true subjective priors, and that inference would be similar to those in the NV-no-information treatment, while distinct from NV-information treatment. Indeed, as we show below, this is precisely what we find. The NV (2007) findings are qualitatively similar to the NV-no-information treatment, while more centered reports in the NV-information treatment results in qualitatively different inference.

Looking at results in the NV-information and NV-no-information treatments, and ignoring beliefs, both replications match the key result of the original study, that conditional on performance, men more than women enter the tournament. Our focus though is on how the elicited beliefs on relative rank move across treatment and how this affects inference.

Focusing first on beliefs, Figure 8 reports the average weight attached to each rank by treatment and gender. Panel A shows the beliefs elicited in the NV-no-information treatment and reveals, as in NV, that men are more confident than women, reporting a substantially greater likelihood of being ranked 1st and a corresponding lower likelihood of being ranked 3rd and 4th.

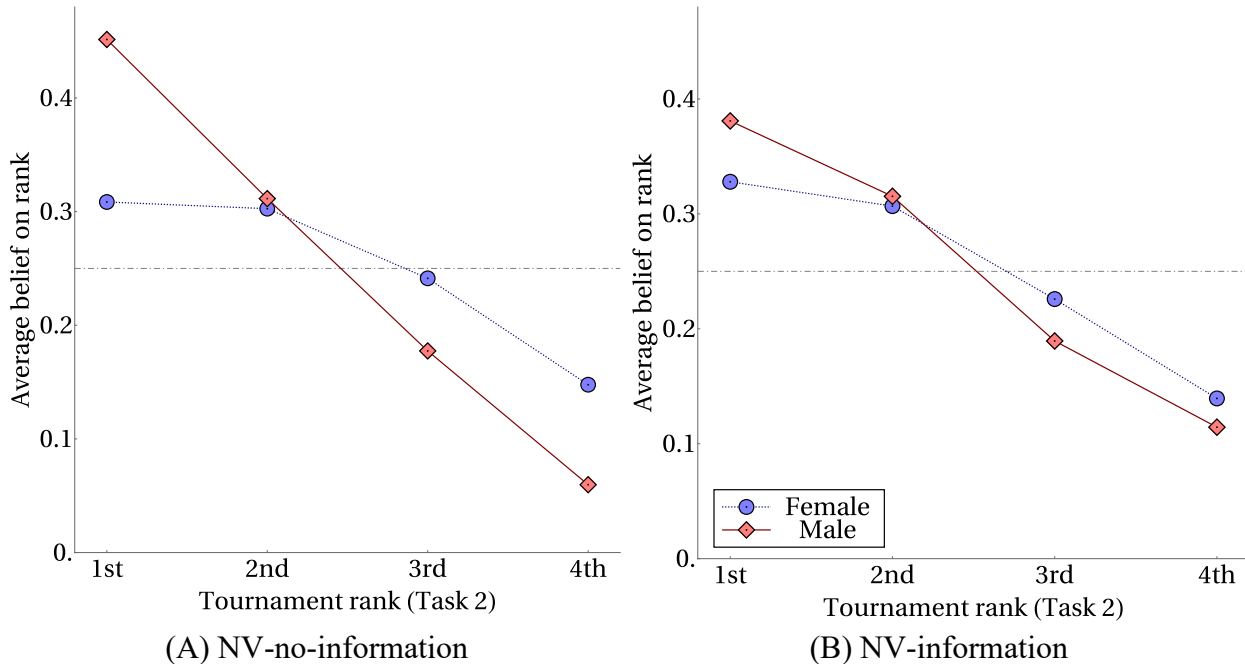


Figure 8. Elicited Likelihood of Performance Rank in Tournament

Our results on eliciting objective priors revealed that information on incentives biased reports on non-centered, but not on centered priors.⁴⁵ With men's beliefs less centered than women's in the NV-no-information treatment we therefore expect that quantitative information on incentives will lead to more center-biased reporting for men than for women. Panel B confirms this expectation in the data. Relative to NV-no-information the reports by men pull-to-center, leading to a much smaller gender gap in confidence for NV-information.

Inferential impact of center bias

Next, we explore the impact of using center-biased belief elicitation for inference. Seeing our NV-no-information sample as capturing the true effect, we can predict and sign the theoretical consequences of center bias in the two inferential exercises from the NV-study. First, when the elicited belief q_i is used as a left-hand-side dependent variable, where we are looking for a difference in means across gender. Second as a right-hand-side independent control variable, where we are looking to draw inference on how the participant's tournament entry decision y_i differs between men and women, after controlling for confidence differences via the elicited beliefs. Each exercise seeks to make inference over a difference between women and men, δ , and our theoretical analysis examines how center-biased reporting might distort inference over δ in the NV-information treatment.

The first key inferential regression in the NV-study uses measured beliefs to examine differences in confidence between men and women:

$$(1) \quad q_i = \mu_q + \delta_q \cdot \text{Female}_i + \epsilon_i,$$

where the estimated gender-gap in confidence $\hat{\delta}_q$ is the focus.⁴⁶ The second inferential exercise examines the tournament-entry decision y_i , conditional on confidence, where the modeled regression is:

$$(2) \quad y_i = \mu_y + \delta_y \cdot \text{Female}_i + \beta_q \cdot q_i + v_i.$$

Here the focus is the estimated gender gap in tournament-entry, $\hat{\delta}_y$, after controlling for the estimated confidence-effect on entry (the $\hat{\beta}_q \cdot q_i$ term).

To predict the inferential distortions we use a simple model of center-bias, where the elicited belief q_i is modeled as a random variable: with probability α the belief is the center-value c (a constant) and with probability $1 - \alpha$ we observe the true belief q_i^* . Using this

⁴⁵ In explaining the substantial change in reports by men, we also find in our elicitation of objective priors a greater center bias for priors further from center (a response that is also consistent with the intuition provided for intensified hedging in Table 1. Using a simple center-biased model, where subject i 's report q_i is given by $q_i = (1 - \alpha) \cdot \pi_i + \alpha \cdot c$, which is an α -weighted average of the centered belief c and her true belief π_i , we estimate significantly greater values of α on priors of 0.2 and 0.8 than on 0.3 and 0.7 ($\hat{\alpha} = 0.289$ vs. $\hat{\alpha} = 0.149$, respectively, $p=0.011$).

⁴⁶ While the actual regressions controls for ability, for simplicity of presentation we focus on the simplest specification to outline the main effect.

model of center-bias, we can predict how inference over $\hat{\delta}_q$ and $\hat{\delta}_y$ will be affected. To facilitate the prediction, we make the following simplifying assumptions: (i) the sample is balanced with N men and N women; and (ii) the econometric errors ϵ_i and ν_i are independent mean zero errors, with finite variance (given by σ_ϵ^2 and σ_ν^2 , respectively).

What then is the effect on inference for center-biased distortions in the beliefs? We here focus on statements of the results, where more detailed derivations are in Online Appendix B.1. When there is no center bias ($\alpha = 0$) the OLS coefficients are unbiased, consistent estimators of the true gender effects δ_q and δ_y ; however, both estimators are biased when measured beliefs are center-biased ($\alpha > 0$).

When beliefs are the dependent variable in (1), center bias moves average beliefs for both men and women to the same point c . As such, the expected *difference* between the two populations is directly attenuated to $(1 - \alpha) \cdot \delta_q$,⁴⁷ leading to an underestimate of the true magnitude for the gender gap under center bias.⁴⁸

When center-biased beliefs are instead used as an independent variable to control for the belief in (2), the asymptotic biases of the OLS estimator are again increasing in the intensity of center bias α (though here in a non-linear manner). However, the bias size and direction also depend on the true population parameters that we are attempting to estimate. In particular, for estimation equation (2) we show that the asymptotic bias is proportional to $\beta_q \cdot \delta_q$, the product of the true marginal effect of confidence on tournament entry and the true gender difference in confidence.⁴⁹ If as in the original NV-study men are more overconfident than women (so $\delta_q < 0$) and confidence has a positive effect on tournament entry ($\beta_q > 0$), then center-biased beliefs are predicted to generate a more-negative estimate for the tournament entry differences between women and men after controlling for beliefs. That is, we predict that $\hat{\delta}_y < \delta_y < 0$. Center-biased beliefs would therefore *overestimate* the size of the gender gap in tournament entry after controlling for beliefs.

Inference from the NV-no-information and NV-information treatments, provided side by side in Table 4, are fully consistent with the predicted effects from center bias. Examining first the case where beliefs are used as the dependent variable in columns (1) and (2), we find that whether or not information is provided on the BSR incentives has a clear qualitative effect on the inference made. The results from the NV-no-information treatment replicate the original NV-finding that men conditional on performance are more confident than women (a 15 percentage-point gap in the believed likelihood of winning the

⁴⁷ The predicted effect of attenuation at rate $1 - \alpha$ would be identical if we were trying to estimate the marginal effect on beliefs for a continuous right-hand-side treatment.

⁴⁸ Where the coefficient is attenuated in proportion to $(1 - \alpha)$, the T -statistic on the $\hat{\delta}_q$ coefficient is more complicated as the variance can be bigger or smaller with the center-biased reports, depending on the location of the center point c relative to the true belief means μ_q and $\mu_q + \delta_q$. However, the variance effect can be bounded, so that the effective t -statistic on the difference in means variable is attenuated by at least $\sqrt{1 - \alpha}$ under the model of distortion.

⁴⁹ In general, when beliefs are on the right-hand-side, asymptotic biases will depend on the covariance between the unobserved belief mismeasurement, the other variables, and the true parameters.

tournament, $p = 0.005$). In contrast, there is no significant gender gap in confidence for NV-information (a 4 percentage-point gap, $p = 0.523$).

Table 4 Gender Differences in Confidence and Tournament Entry: NV replication Results with and without Information.

	DEPENDENT VARIABLE (CF. TABLE V IN NV 2007)		INDEPENDENT VARIABLE (CF. TABLE II AND VI IN NV 2007)			
	Belief on 1st rank (OLS)		Tournament entry (Probit)			
	No-Inform.	Information	No-Information		Information	
	(1)	(2)	(3)	(4)	(5)	(6)
Female	-0.148 (0.051)	-0.038 (0.059)	-0.258 (0.080)	-0.146 (0.115)	-0.357 (0.118)	-0.382 (0.124)
Tournament	0.008 (0.004)	0.018 (0.005)	0.022 (0.012)	0.017 (0.014)	0.006 (0.011)	-0.011 (0.014)
Tournament– piece rate	0.017 (0.008)	-0.021 (0.010)	0.001 (0.017)	-0.015 (0.021)	-0.023 (0.021)	-0.004 (0.022)
Constant	0.305 (0.097)	0.059 (0.101)				
Belief weight on 1st rank				1.275 (0.432)		0.994 (0.329)
N	74	68	74	74	68	68
R ² /adj. R ²	0.273	0.187	0.157	0.303	0.093	0.208

Note: Columns 1-2: Tobit regressions yield the same qualitative results. Columns 3-6: Marginal effects (estimated constants omitted). Regressions including the average reported rank instead of the reported weight on rank 1 yield the same qualitative results.

Next, exploring the consequences of using center-biased beliefs as an independent variable we examine the gender gap in tournament entry after controlling for confidence. Again, the results of the NV-no-information study mirror those of the original. Conditional on performance, men more than women enter the tournament (column 3), with the gender gap in tournament entry being partly explained by the gender gap in confidence (column 4). Participants that think they will win the tournament are more likely to enter it, and so controlling for beliefs substantially reduces the gender gap in entry. The results for the same estimation procedure in our NV-information sample (columns 5 and 6) lead to a very different conclusion.⁵⁰ While conditional on performance men are more likely to select the

⁵⁰ Fully replicating the NV (2007) design with the original coarse belief elicitation a recent study by Recalde and Vesterlund (2022), with 263 participants at the University of Pittsburgh, finds that women were 25.2 percent less likely than men to report a belief that they were first in their group (corresponding to Table 4 columns 1 and 2). Further an initial gender gap in tournament entry of 27.8 percent was shown to decrease to 18.9 percent when controlling for the coarse belief measure of being ranked first (corresponding to Table 4 columns 3 and 4 versus 5 and 6). All gender differences were significant at the 0.001 level.

tournament (column 5), controlling for beliefs does not reduce the gap in tournament entry (column 6).⁵¹

Our experimental results match what our theory predicts, that center-biased reporting in NV-information can lead to the false inference that there is no gender gap in confidence and that confidence plays no role in explaining the gender gap in competition.⁵²

Further, the results from the NV-information treatment mirror the comparative statics when we simulate center-biased beliefs using the data from the NV-no-information treatment (distorting each belief q_i in the data with probability α to the center-value c). The simulations also indicate a smaller gender gap in confidence and an exaggerated assessment of the gender gap in tournament entry after controlling for the beliefs. By contrast, attempts to reconstruct the NV-no-information results using the NV-information data are not successful. Once the center-bias masks the gender gap in confidence, it is not easily uncovered (see Appendix B.2. for details).⁵³

Our Information-No-Information replication of the NV-study confirms that information on incentives leads to centered reporting, here over subjective beliefs. Moreover, the direct comparison demonstrates the substantial impact of using biased beliefs for inference, with center-biased reporting causing us to underestimate the gender gap in confidence and overestimate the gender gap in preferences for competition.

5.2. Implication for belief elicitation

Our NV-replication demonstrates the serious inferential impact of using an elicitation that fails a very weak condition for behavioral incentive compatibility. What is less clear is how scholars should respond to the finding. If incentives lead to distorted reports under this state-of-the-art belief elicitation, then how do we elicit beliefs?

Needless to say, it is problematic to use any elicitation where information on the incentives underlying that mechanism increases the rate of false reports. But perhaps information only distorts reports under BSR? To demonstrate the diagnostic merit of our Information-No-Information comparison we apply it to an examination of the quadratic

⁵¹ To secure no prior exposure to similar studies we conducted the study with 147 students from UCSB. Adding an ‘other’ option to the demographic question on gender we have 5 students in the NV-no-information treatment who do not identify as male or female. Mirroring the NV-analysis we examine responses only for the 142 participants who identify as male or female. Including all participants, and assessing the difference between the either male/not-male or not-female/female yields the same effects for both Information and No Information. While not-female/female yields the same qualitative effects as those shown in Table columns (1) and (4), the not-male coefficient is significant in the equivalent of column (4) ($p = 0.041$).

⁵² Considering the broader impact of the NV study this false inference could well impact policies aimed to address gender differences in advancement (e.g., the study was in a Submission to the Senate’s Employment, Workplace Relations and Education Legislation Committee’s Inquiry into the Workplace Relations Amendment (WorkChoices) Bill 2005 in Australia: http://www.hreoc.gov.au/legal/submissions/workplace_relations_amendment_2005.html.

⁵³ Assessment of potential inferential bias in prior work therefore calls for Information-No-Information replications of the study of interest.

scoring rule (QSR).⁵⁴ The results mirror those from the BSR. Information on the QSR incentives increases the rate of false reports—persistently over time and particularly for non-centered priors (39.2 percent false reports in *QSR-information* versus 25.6 percent in *QSR-no-information*; see Figure A.5 in the Online Appendix). The QSR’s failure reflects the more widely accepted idea that this elicitation is not incentive compatible for risk averse agents. However, it also corroborates the diagnostic power of our Information-No-Information treatment comparison in determining whether an elicitation violates weak conditions over behavioral incentive compatibility.⁵⁵

When considering improved elicitations, the finding that information on the incentives distort reports, tempts the somewhat perverse option of reducing the information provided. One option here is to simply rely on the (truthful) qualitative statement that accurate reporting maximizes the chance of winning. Certainly, the data from our No-Information treatment where only this statement is given, shows that this is the better option in terms of the accuracy of collected belief data. However, advocating for what amounts to a fully black-box mechanism from the point of participants is jarring to the general philosophy of incentivized decision making.

A less extreme option is to also add a description of the mechanism’s implementation rule, without providing precise quantitative details on the incentives. As seen in Table 3 this procedure has been used in recent implementations.⁵⁶ We pursue this approach in what we refer to as our *Description treatment*. We augment the limited information in No-Information (a statement on dominance), with the short non-quantitative description, used in the Information treatment, of how prize realizations are determined—but with no other information provided.⁵⁷ Mathematically inclined participants are thus informed on the

⁵⁴ The experiments were conducted online using 60 University of Pittsburgh participants in each treatment using our online protocol.

⁵⁵ While we cannot rule out that the center-biased reporting under the QSR results from risk-aversion, we see no evidence that it does. Individual false-report rates and center-biased reporting are not significantly correlated with attitudes towards risk (whether measured as the certainty equivalent of a gamble, or as a probability equivalent switch point from gamble to a certain option, all $p > 0.284$).

⁵⁶ This approach is frequently used in other mechanisms. For example, consider the non-technical description of how a second-price bidding rule works (and equivalently, how strategy methods like the BDM function), or how a complicated matching algorithm like top-trading cycles would be described to parents providing school-choice rankings. Our results (with a fixed mechanism) dovetail with Holt and Smith (2016) who find evidence across mechanisms for the superiority of a BDM-based crossover elicitation. Similar to our Description treatment, their crossover mechanism does not spell out the marginal effects on the probability of winning, focusing on the qualitative compatibility. In comparison, their QSR elicitation uses a table to make clear the marginal effects on the monetary prize.

⁵⁷ Participants received no quantitative information on the incentives prior or after submitting a report (so no information on the lottery pair at each choice, and no feedback after the round). Information on the incentives were only provided at the start of the experiment, augmenting the instructions for the No-Information treatment with a compact description of the payment rule. Specifically, we make use of two uniform draws, and tell participants that they will win the \$8 if: (i) The event happens and their report is greater than the smaller of the two random draws; or (ii) the event does not happen and their reported belief is lower than the larger of the two random draws. As shown in Wilson and Vespa (2017) this rule is payoff equivalent to the BSR. The experiment was conducted in person at PEEL with 60 unique participants.

mechanism's quantitative incentives, while the less mathematically inclined learn that a concrete procedure is used to map reported beliefs into final earnings.

Results reported at the bottom of Table 2 reveal that the Description treatment largely mirrors that of the No-Information treatment, with a moderate rate of false reports (24.5 percent). This rate is not significantly different from the false-report rate in No Information ($p=0.610$) but is significantly lower than the Information treatment ($p=0.004$). As in the No-Information treatment, there is no evidence that false reports pull-to-center ($p=0.564$), though false-report rates are slightly, but significantly, higher for non-centered than for centered priors ($p=0.019$).⁵⁸

As such, the complete ambiguity over the elicitation's incentives in No Information can be relaxed without severely damaging reports. However, given the distinctly different reporting behavior when participants receive detailed information on the quantitative incentives (Information treatment), it is unlikely that participants comprehend the offered incentives in Description.⁵⁹

While we may be able to limit false reporting by presenting the barest possible information on incentives, this is hardly what we have in mind when designing incentive compatible elicitation. The intent is for the incentives to drive participants toward truth telling rather than away from it.

We see our consistent finding that information on the BSR incentives increases false reports as a violation of a very weak condition for behavioral incentive compatibility. A related but more-direct assessment of the incentives evaluates another weak condition for behavioral incentive compatibility, that most participants select the outcome assumed to be maximizing under the mechanism.

Fixing again the BSR incentives as the example, we explore this condition by conducting an *Incentives-only treatment*, where participants are given a choice over the set of lottery pairs underlying the BSR, stripped of the elicitation framing. The set of 11 lotteries mirror those shown in Table 1, with each pair of tickets consisting of a red- and a blue-lottery ticket. One lottery ticket in the participant's chosen pair is implemented for payment, where participants are informed that the ticket-color counting for payment is

⁵⁸ Figures of the rates of false reports by prior and by period are shown in Figure A.6 in the online appendix. Intriguingly, while there is no time trend in our Information, No-Information or RCL treatments, the rate of false reports decreases over time in the Description treatment ($p=0.077$). This decrease in false reports is consistent with the information on incentives becoming less salient in the Description treatment, where participants only learn the incentives up front. This salience argument can be mirrored in the Feedback treatment where false reports increase as participants get repeated feedback on the incentives.

⁵⁹ Looking at the response to survey questions and coding these as agreeing to the statement on understanding how payoffs were calculated, how a stated belief affected pay, and whether they truthfully reported we find that the Description and the No-Information treatments are statistically indistinguishable from one another in participants' self-reported understanding of the mechanism (77 vs. 72 percent, $p=0.532$, χ^2 -test) but that there are differences in understanding how beliefs affected pay (83 vs. 70 percent, $p=0.084$), and indications of differences in self-reported inclination to report truthfully (75 vs. 85 percent, $p=0.171$). Reporting an understanding of how beliefs affected pay may result from understanding that truthful reporting maximized the chance of winning the prize.

determined by a random draw with a predetermined chance (20 or 30 percent) on the red-ticket lottery counting.⁶⁰ Truthful revelation under the BSR relies on participants perceiving one unique lottery pair as the maximizer for each possible probability on the elicited event.

In Figure 9 we illustrate the lottery choices in the Incentives-only treatment. On the left of each panel we present the proportion choosing the theorized maximizer, while on the right we report the deviations (here identifying the chosen lottery with the corresponding prior belief under the BSR rather than its lottery label in the experiment). We see in violation of our weak condition for behavioral incentive compatibility that the majority of participants fail to select the lottery thought to be maximizing. In fact, the distribution shows wider dispersion than observed under the comparable prior elicitation (cf. Figure 5).⁶¹ Common to both sets of choices though is a tendency to select lottery pairs that have similar chances of winning, reflecting the hedging-motives alluded to in the introduction, and consistent with the center-biased reporting documented in our Information-No-Information comparisons.⁶²

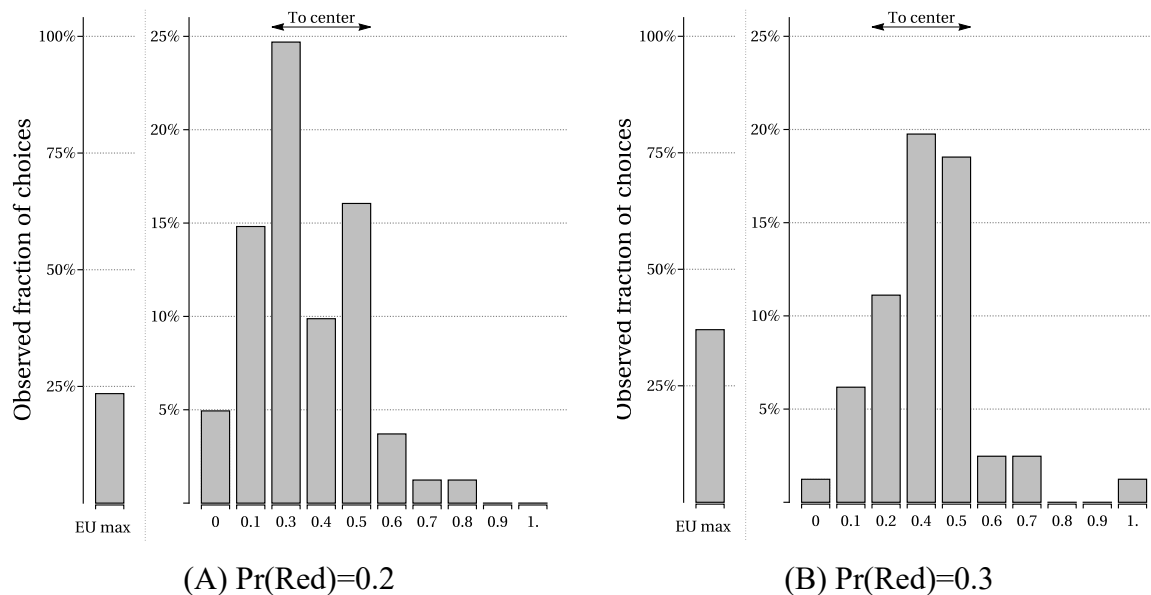


Figure 9: Chosen outcomes under Incentives-Only Treatment

⁶⁰ The exact choices are made over lottery pairs A through K, with the chance of winning for lottery A being 100 percent on the red lottery ticket and 0 percent on the blue, for B it is 99 percent on red and 19 percent on blue, etc. This short study was administered to 162 participants as a module at the end of another PEEL study examining public-goods provision, similar to the implementation of Fischbacher and Föllmi-Heusi (2013). One of the two decisions was selected for payment, with an \$8 prize in the lottery.

⁶¹ We inform participants when eliciting beliefs that accurate reporting maximizes the chance of winning, but in contrast participants in the Incentives-only treatment are given no guidance on what the 'right' choice is.

⁶² For the Incentives-only treatment 50.6 percent (38.3 percent) of the participants in the $\text{Pr}(\text{Red})=0.20$ (0.3) condition move towards the center, the illustrated choice range in Figure 9. In contrast the proportion moving toward the nearest extreme is 19.8 percent (18.5 percent).

What then are the implications of these findings for belief elicitation more generally? As the state-of-the-art elicitation, the BSR conveniently allows us to truthfully state to participants that their chances of winning are maximized by truthful reporting. However, the evidence for violations of two weak conditions—that provision of quantitative information on the incentives increases false reporting, along with the finding that participants when presented with the bare incentives mostly fail to select the presumed maximizing outcome—shows that the elicitation is not behaviorally incentive compatible. Adding to this concerning evidence is the clear impact that the resulting center-biased reporting has on inference.

In pursuing improved elicitations, we need to be cognizant that we are designing mechanisms for behavioral agents. In this respect, our findings and proposed tests for behavioral incentive compatibility relate to Li's (2017) concepts of obvious dominance and obvious strategy proofness. Both our work and Li (2017) stress the importance of considering cognitive limitations (in addition to a broader set of preferences) when designing incentive compatible mechanisms. However, while Li (2017) provides a theoretical criterion of a mechanism's incentive compatibility for a class of cognitively limited agents, our work stresses the importance of, and provides means to, testing whether a theoretically incentive compatible mechanism is behaviorally incentive compatible in an empirical sense. As in the BSR, relatively weak-seeming theoretical assumptions permit the design of fully-separating mechanisms, to measure beliefs at arbitrary precision. But such precision may well be costly—where we need to empirically test that the assumptions put in place hold, and that behavioral agents actually perceive truthful revelation as beneficial.

Our study has proposed weak conditions for behaviorally incentive compatible elicitations and provided diagnostic tools for checking them. The hope is that new elicitations will be assessed against and succeed in passing these standards. Given the challenges associated with this task though, we caution that it may be time to question whether it is reasonable to assume that participants in our studies hold exact probabilistic beliefs, let alone our ability to use monetary incentives to elicit such beliefs at arbitrary precision. Instead of taking our results as a call for the development of mechanisms that are incentive compatible for an ever-more-general class of decision maker, we might instead ask whether the necessary economic inferences could be drawn with less-precise measurements, where the incentives for truthful reporting can be simpler and starker.⁶³ For example, in discrete settings it may be sufficient to elicit the event the participants deem most likely and incentivize the elicitation by offering compensation only in the event that

⁶³ Ex post corrections of beliefs à la Offerman et al (2009) similarly relies on individuals holding exact beliefs. If precise beliefs are a requirement, then dynamic mechanisms that elicit the same belief through multiple (adaptive) coarse elicitations may be used to make incentives appear stark (see Schmidt and Zankiewicz, 2016). With starker incentives we mean mechanisms that provide steeper marginal incentives, especially in the neighborhood of the theorized maximizer.

the report is correct.⁶⁴ In continuous settings, the same can be achieved by paying participants if the true population outcome falls within some bounds around their guess.⁶⁵ Alternatively, it may be sufficient to determine whether a belief lies within a certain *fixed* interval. This allows for deviations between the potential intervals to come at a higher perceived cost and may still provide the information necessary for inference.⁶⁶ For example, suppose that in understanding individual behavior we wish to elicit the belief that an opponent will select action A or B , and that the individual's predicted behavior theoretically depends on the belief on A exceeding a 30 percent cutoff. Rather than eliciting the precise belief that action A is chosen, it may secure more reliable and truthful reporting to instead focus the elicitation on whether or not the belief on A exceeds the theoretical cutoff. If elicited beliefs are collected primarily as controls or for auxiliary tests of a behavioral mechanic, inference may be improved with starker incentives over coarser elicitations.

While there are many paths to improve belief elicitation, we propose two simple assessments: that information on the incentives increases truthful reporting; and that most participants when given a choice over the pure set of incentives select the theorized maximizer. In demonstrating the very substantial inferential consequences from using biased elicitations, our results serve as a call for elicitations to be incentive compatible both theoretically and behaviorally, but also as a strong caution against elicitations that rely on incentives that decrease truthful reporting.

⁶⁴ See, for example, Bhatt and Camerer, 2005, Hurley and Shogren, 2005, Niederle and Vesterlund, 2007, Vanberg, 2008, Blanco et al., 2010, Dargnies, 2012, Di Tella et al. 2015, LeCoq et al., 2015, Toussaert, 2018, Bordalo et al. 2019, Cantoni et al., 2019, Wilcox and Feltovich, 2000, and Huffman et al., 2019

⁶⁵ For example, Charness and Dufwenberg, 2006, Abeler et al, 2019, Danz et al., 2018. As Abeler et al (2019) point out, “[t]his mechanism is very simple and easier to explain and understand than proper scoring rules. It elicits in an incentive-compatible way the mode (or more precisely, the mid-point of the [x]-percentage point interval with the highest likelihood) of a subject’s distribution of estimates.”

⁶⁶ Asking participants to select one of a several fixed ranges of probabilities acknowledges both the ambiguity associated with providing an exact probabilistic belief (Manski, 2004) and makes it possible to provide starker incentives for doing so.

REFERENCES

- Abdellaoui, Mohammed, Driouchi, Ahmed, and L'Haridon, Olivier, "Risk aversion elicitation: reconciling tractability and bias minimization," *Theory and Decision* 71, 1 (2011), pp. 63–80.
- Abeler, Johannes, Nosenzo, Danielle, Raymond, Collin, "Preferences for Truth-Telling," *Econometrica*, 87, 4, (2019), pp. 1–118.
- Ahrens, Steffen, and Ciril Bosch-Rosa. "The Motivated Beliefs of Investors Under Limited Liability." Working paper, 2019.
- Aksoy, Billur, Ian Chadd, and Boon Han Koh. "Hidden Identity and Social Preferences: Evidence from Sexual Minorities." Working paper, 2021.
- Aoyagi, Masaki, Guillaume R Fréchette, and Sevgi Yuksel. 2021. "Beliefs in Repeated Games." Working paper, 2021.
- Avoyan, Ala and Schotter, Andrew, "Attention in Games: An Experimental Study" (2020), *European Economic Review*, 124, 103410 (2020), pp. 1–28.
- Babcock, Linda, Recalde, Maria P, Vesterlund, Lise, and Weingart, Laurie, "Gender differences in accepting and receiving requests for tasks with low promotability," *American Economic Review* 107, 3 (2017), pp. 714–47.
- Benjamin, Daniel, "Errors in Probabilistic Reasoning and Judgment Biases," in Bernheim Douglas, DellaVigna, Stefano, Laibson, David, eds., *Handbook of Behavioral Economics: Applications and Foundations 1*, Volume 2, North-Holland, (2019), pp. 69–186.
- Bhatt, M, Camerer, CF. "Self-referential thinking and equilibrium as states of mind in games: fMRI evidence," *Games and Economic Behavior* 52, 2, (2005), pp.424–459.
- Blanco, M, Engelmann, D, Koch, AK, Normann H-T, "Belief elicitation in experiments: is there a hedging problem?" *Experimental Economics*, 13, 4, (2010), pp.412–438.
- Bordalo, Pedro, Coffman, Katherine, Gennaioli, Nicola, and Shleifer Andrei, "Beliefs about Gender," *American Economic Review* 109, 3 (2019), pp.739–773.
- Bruner, David M, "Changing the probability versus changing the reward," *Experimental Economics* 12, 4 (2009), pp. 367–385.
- Burfurd, Ingrid and Wilkening, Tom, "Experimental guidance for eliciting beliefs with the Stochastic Becker–DeGroot–Marschak mechanism." *Journal of the Economic Science Association* 4, 1, (2018), pp. 15–28.
- Cantoni, D., Yang, D.Y., Yuchtman, N, Zhang, YJ., "Protests as Strategic Games: Experimental Evidence from Hong Kong's Antiauthoritarian Movement," *Quarterly Journal of Economics*, 134, 2, (2019), pp. 1021–1077.
- Cason, Timothy N. and Plott, Charles R. "Misconceptions and Game Form Recognition: Challenges to Theories of Revealed Preference and Framing," *Journal of Political Economy*, 122 (2014), pp. 1235–1270.

- Cason, Timothy N., Sharma, Tridib, and Vadovič, Radovan. “Correlated beliefs: Predicting outcomes in 2×2 games.” *Games and Economic Behavior*, 122 (2020), pp. 256–276.
- Castillo, Marco E, Philip J Cross, and Mikhail Freer. 2019. “Nonparametric Utility Theory in Strategic Settings: Revealing Preferences and Beliefs from Proposal–Response Games.” *Games and Economic Behavior* 115 (2019), pp. 60–82.
- Charness, Gary, Dufwenberg, Martin “Promises and Partnership,” *Econometrica*. 74, 6, (2006), pp. 1579–1601.
- Charness, Gary, Oprea, Ryan, and Yuksel, Sevgi. “How do people choose between biased information sources? Evidence from a laboratory experiment.” *Journal of the European Economic Association*, 19:3, (2021), pp. 1656-91.
- Chen, Yan, and YingHua He. n.d. “Information Acquisition and Provision in School Choice: A Theoretical Investigation.” forthcoming *Economic Theory*, (2021), 10.1007/s00199-021-01376-3
- Choi, Syngjoo, Byung-Yeon Kim, Jungmin Lee, and Sokbae Lee. “Institutions, Competitiveness and Cognitive Ability.” Working paper, 2020.
- Colzani, Paola, and Luis Santos-Pinto. “Does Overconfidence Lead to Bargaining Failures?” Working paper, 2020
- Corazzini, Luca, Stefano Galavotti, and Paola Valbonesi. “An Experimental Study on Sequential Auctions with Privately Known Capacities.” *Games and Economic Behavior* 117 (2019), pp. 289–315.
- Danz, David, Madarász, Kristóf, Wang, Stephanie, “The Biases of Others: Projection Equilibrium in an Agency Setting,” *CEPR Discussion Paper*, DP12867, (2018).
- Danz, David, Gupta, Neeraja, Lepper, Marissa, Vesterlund, Lise, and Winichakul, K. Pun, “Going virtual: A step-by-step guide to taking the in-person experimental lab online,” working paper, 2021.
- Danz, David, Vesterlund, Lise and Wilson, Alistair “Data and Code for : Belief Elicitation and Behavioral Incentive Compatibility.” *American Economic Association* [publisher], Inter-university Consortium for Political and Social Research [distributor]. <https://doi.org/10.3886/E157161V1>.
- Dargnies, Marie-Pierre, Hakimov, Rustamdjan, and Kübler, Dorothea, “Self-Confidence and Unraveling in Matching Markets,” *Management Science* 65, 12 (2019), pp. 5603–18.
- Dargnies, Marie-Pierre, “Men Too Sometimes Shy Away from Competition: The Case of Team Competition,” *Management Science* 58, 11, (2012), pp. 1982–2000.
- Dean, Mark, and Ortoleva, Pietro, “Allais, Ellsberg, and Preferences for Hedging,” *Theoretical Economics* 12 (2017):377-424.
- DeQuidt, Jonathan, Vesterlund, Lise, and Wilson, Alistair, “Experimenter Demand Effects,” *Handbook of Research Methods and Applications in Experimental Economics* (Edward Elgar Pub., edited by Arthur Schram and Aljaz Ule).

- Dianat, Ahrash, Echenique, Federico, and Yariv, Leeat, “Statistical Discrimination and Affirmative Action in the Lab,” *CEPR Discussion Paper No. DP12915*, (2019).
- Di Tella, Rafael, Perez-Truglia Ricardo, Babino Andres, Sigman Mariano, “Conveniently Upset: Avoiding Altruism by Distorting Beliefs about Others’ Altruism,” *American Economic Review* 105, 11 (2015), pp. 3416–3442.
- Drerup, Tilman, Enke, Benjamin, and Von Gaudecker, Hans-Martin, “The precision of subjective data and the explanatory power of economic models,” *Journal of Econometrics* 200, 2 (2017), pp. 378–389.
- Dustan, Andrew, Kristine Koutout, and Greg Leo. “Second-Order Beliefs and Gender.” Working paper, (2020)
- Enke, Benjamin, and Graeber, Thomas, “Cognitive Uncertainty”. *NBER Working Paper No. 26518*, (2019).
- Enke, Benjamin, Schwerter, Frederik, and Zimmermann, Florian, “Associative memory and Belief Formation,” *NBER Working Paper No. 26664*, (2020).
- Erkal, Nisvan, Lata Gangadharan, and Boon Han Koh. 2020. “Replication: Belief Elicitation with Quadratic and Binarized Scoring Rules.” *Journal of Economic Psychology* 81: 102315.
- . “By Chance or by Choice? Biased Attribution of Others’ Outcomes When Social Preferences Matter.” Working paper, (2021).
- . “Attribution biases in leadership: Is it effort or luck?” Working paper, (2019).
- Filippin, Antonio, and Marco Mantovani. “Risk Aversion and Information Aggregation in Asset Markets.” Working paper (2019).
- Fischbacher, Urs, “z-Tree: Zurich toolbox for ready-made economic experiments,” *Experimental Economics* 10, 2 (2007), pp. 171–178.
- and Franziska Föllmi-Heusi, “Lies in Disguise—An Experimental Study on Cheating.” *Journal of the European Economic Association*, 11(3), (2013), pp 525–547,
- Frederick, Shane, “Cognitive reflection and decision making,” *Journal of Economic Perspectives* 19, 4 (2005), pp. 25–42.
- Gächter, Simon and Renner, Elke, “The effects of (incentivized) belief elicitation in public goods experiments,” *Experimental Economics* 13, 3 (2010), pp. 364–377.
- Graeber, Thomas. “Inattentive Inference.” Working paper, (2021).
- Hao, Li and Houser, Daniel, “Belief Elicitation in the Presence of Naïve Respondents: An Experimental Study,” *Journal of Risk and Uncertainty* 44, 2 (2012), 161–80.
- Harrison, Glenn W and Phillips, Richard D, “Subjective beliefs and statistical forecasts of financial risks: The chief risk officer project,” in *Contemporary Challenges in Risk Management* (Springer, 2014), pp. 163–202.

- Healy, Paul J, “Epistemic Game Theory Experiments: Utility Elicitation and Irrational Play” mimeo, (2017).
- Healy, Paul J, “Explaining the BDM—Or any Random Binary Choice Elicitation Mechanism—To Subjects”, mimeo, (2018).
- Hillenbrand, Adrian and Schmelzer, André, “Beyond information: Disclosure, distracted attention, and investor behavior,” *Journal of Behavioral and Experimental Finance* 16 (2017), pp. 14–21.
- Holt, Charles A. and Smith, Angela M., “Belief Elicitation with a Synchronized Lottery Choice Menu That Is Invariant to Risk Attitudes,” *American Economic Journal: Microeconomics* 8, 1 (2016), pp. 110–39.
- Hossain, Tanjim and Okui, Ryo, “The binarized scoring rule,” *Review of Economic Studies* 80, 3 (2013), pp. 984–1001.
- , “Belief formation under signal correlation.” Working paper, (2019).
- Huffman, David, Raymond, Collin, Shvets, Julia, “Persistent Overconfidence and Biased Memory: Evidence from Managers,” mimeo, (2019).
- Hurley, TM, Shogren, JF. “An Experimental Comparison of Induced and Elicited Beliefs,” *Journal of Risk and Uncertainty*, 30, 2, (2005), pp. 169–188.
- Koutout, Kristine. 2020. “Gendered Beliefs and the Job Application Decision: Evidence from a Large-Scale Field and Lab Experiment.”
- LeCoq, C, Tremewan, J, Wagner, AK., “On the effects of group identity in strategic environments,” *European Economic Review*, 76, (2015), pp. 239–252.
- Li, Shengwu, “Obviously strategy-proof mechanisms.” *American Economic Review* 107, 11, (2017), pp. 3257–87.
- Manski, Charles F. "Measuring expectations." *Econometrica* 72.5 (2004): 1329–1376.
- Masiliūnas, Aidas, “Overcoming coordination failure in a critical mass game: strategic motives and action disclosure,” *Journal of Economic Behavior & Organization* 139 (2017), pp. 214–251.
- Meloso, Debrah, Nunnari, Salvatore, and Ottaviani, Marco, “Looking into crystal balls: a laboratory experiment on reputational cheap talk” Working paper, (2020).
- Nelson, Robert G and Bessler, David A, “Subjective probabilities and scoring rules: Experimental evidence,” *American Journal of Agricultural Economics* 71, 2 (1989), pp. 363–69.
- Niederle, Muriel and Vesterlund, Lise, “Do women shy away from competition? Do men compete too much?” *Quarterly Journal of Economics* 122, 3 (2007), pp.1067–1101.
- Offerman, Theo, Sonnemans, Joep, Van De Kuilen, Gijs, and Wakker, Peter P, “A Truth Serum for Non-Bayesians: Correcting Proper Scoring Rules for Risk Attitudes,” *The Review of Economic Studies*, 76, 4 (2009), pp.1461–89.

- Oprea, Ryan, and Sevgi Yuksel. "Social Exchange of Motivated Beliefs." *Journal of the European Economic Association*. Forthcoming, (2021). DOI: 10.1093/jeea/jvab035/6368096
- Palfrey, Thomas R and Wang, Stephanie W, "On eliciting beliefs in strategic games," *Journal of Economic Behavior & Organization* 71, 2 (2009), pp. 98–109.
- Rafkin, Charlie, Advik Shreekumar, and Pierre-Luc Vautrey. "When Guidance Changes: Government Stances and Public Beliefs." *Journal of Public Economics* 196, (2021).
- Recalde, Maria, and Lise Vesterlund, "Gender Differences," (2022), University of Pittsburgh working paper
- Renes, Sander, and Bauke Visser. "Markets Assessing Decision Makers and Decision Makers Impressing Markets: A Lab Experiment." Working paper, (2019).
- Roth, Alvin E and Malouf, Michael W, "Game-theoretic models and the role of information in bargaining," *Psychological Review* 86, 6 (1979), pp. 574.
- Schlag, Karl H and van der Weele, Joël J, "A Method to Elicit Beliefs as Most Likely Intervals," *Judgment and Decision Making* 10, 5 (2015), 456–68.
- Schmidt, Tobias and Zankiewicz, Christian, "Binary Choice Belief Elicitation: An Adaptively Optimal Design" (2016).
- Schotter, Andrew and Trevino, Isabel, "Belief Elicitation in the Laboratory," *Annual Review of Economics* 6,1 (2014) pp. 103–28,
- Sonsino, Doron, Yaron Lahav, and Amir Levkowitz. 2020. "The Conflicting Links Between Forecast-Confidence and Trading Propensity." *Journal of Behavioral Finance*, (2020) pp. 1–18.
- Toussaert S. "Eliciting Temptation and Self-Control Through Menu Choices: A Lab Experiment," *Econometrica*, 86, 3, (2018), pp. 859–889.
- Trautmann, Stefan T and van de Kuilen, Gijs, "Belief elicitation: A horse race among truth serums," *Economic Journal* 125, 589 (2015), pp. 2116–35.
- Vanberg, C. "Why Do People Keep Their Promises? An Experimental Test of Two Explanations," *Econometrica*, 76, 6, (2008), pp. 1467–1480.
- Wang, Stephanie W, "Incentive effects: The case of belief elicitation from individuals in groups," *Economics Letters* 111, 1 (2011), pp. 30–33.
- Wilcox NT, Feltovich N. "Thinking Like a Game Theorist: Comment," *Mimeo, University of Houston*, (2000).
- Wilson, Alistair J. and Vespa, Emanuel, "Paired-uniform scoring: Implementing a binarized scoring rule with non-mathematical language" (2018).
- Zimpelmann, Christian. 2021. "Stock Market Beliefs and Portfolio Choice in the General Population." Working paper, (2021).