

# Priors and Payoffs in Confidence Judgments

Shannon M. Locke<sup>1,3\*</sup>, Elon Gaffin-Cahn<sup>1\*</sup>, Nadia Hosseinizaveh<sup>1</sup>, Pascal Mamassian<sup>3</sup>,  
Michael S. Landy<sup>1,2</sup>

**1** Dept. of Psychology, New York University, New York, NY, United States

**2** Center for Neural Science, New York University, New York, NY, United States

**3** Laboratoire des Systèmes Perceptifs, Département d'Études Cognitives, École Normale Supérieure, PSL University, CNRS, 75005 Paris, France

\* indicates authors who contributed equally to this manuscript.

Email of corresponding author: shannon.m.locke@nyu.edu

## 1 Abstract

Priors and payoffs are known to change perceptual decision-making, but little is understood about how they influence confidence judgments. Human observers performed an orientation-discrimination task with varied priors and payoffs. We investigated the subsequent placement of discrimination and confidence criteria by comparing behavior to several plausible Signal Detection Theory models. A normative account of behavior uses optimal discrimination criteria. Optimal confidence criteria are yoked to the accuracy-maximizing criterion (i.e., are not affected by payoffs). Additionally, in a normative account, the criterion shifts predicted for asymmetric payoffs and priors should sum when both are varied. We found that observers were conservative in discrimination-criterion placement and that criterion shifts due to priors and payoffs did not sum. For confidence judgments, observers exhibited one of two sub-optimal behaviors. One subset of observers used fixed confidence criteria independent of priors and payoffs. The other group of observers always shifted their confidence criteria with the gains-maximizing discrimination criterion. Such metacognitive mistakes about one's perceptual choices could have negative consequences outside the laboratory setting.

**Keywords:** *decision-making, metacognition, confidence, Signal Detection Theory.*

## 2 Introduction

In making a perceptual decision, it is wise to consider information beyond the available sensory evidence. To maximize expected gains, one should consider both the baseline probability of each possible world state, i.e. *priors*, as well as the associated risks and rewards for choosing or not choosing each response alternative, i.e., *payoffs*. In the Signal Detection Theory (SDT) framework, priors and payoffs alter the threshold amount of evidence required to choose one alternative versus another, that is, a shift in the *criterion* for reporting option “A” versus option “B” in a binary task. For example, a radiologist may be trying to detect a tumor from an x-ray. The radiologist should be more likely to report a positive result for a suspicious shadow if the patient’s file indicates they are a smoker, as this means they have a higher prior probability of cancer. Similarly, the high cost of waiting to treat the cancer should also bias the radiologist towards declaring a positive result. In both real and laboratory environments, observers have been found to factor in priors and payoffs when setting the decision criterion (Maddox and Bohil, 1998, 2000; Maddox and Dodd, 2001; Wolfe et al., 2005; Ackermann and Landy, 2015; Horowitz, 2017), with some caveats we will discuss shortly.

Decisions about the state of the world (cancer or not cancer, cat or dog, clockwise or counter-clockwise of vertical) are based on the stimulus alone and are classified as stimulus-conditioned responses or *Type 1* decisions in the literature. These differ from *Type 2* decisions (or response-conditioned responses), which are judgments about the correctness of Type 1 decisions (Clarke et al., 1959; Mamassian, 2016). In layman’s terms, Type 2 responses are the observer’s confidence about a decision they’ve made, which are often operationalized in binary decision-making experiments as a subjective estimate of the probability the Type 1 response was correct (Pouget et al., 2016). Confidence plays a broad role in guiding behavior, subsequent decision-making, and learning in a multitude of scenarios for both humans and animals (Metcalf and Shimamura, 1996; Smith et al., 2003; Beran et al., 2012).

How does an ideal-observer radiologist modify confidence judgments in response to varying priors or payoffs? Intuitively, a radiologist should be more confident in a positive diagnosis when the patient is a smoker, given the prior scientific literature on the health risks

of smoking that the radiologist has read. Additional confirmatory information should boost 46  
confidence in that positive diagnosis, and contrary evidence should reduce confidence, be- 47  
cause priors (smoker or non-smoker) and sensory evidence (cancerous-looking shadow) are 48  
both informative about the likelihood over possible world states. However, this is not the 49  
case for payoffs. Incentivizing the different responses with rewards or costs does not change 50  
the uncertainty about the world state. The radiologist should not be more or less sure of a 51  
cancer diagnosis if the type of cancer would be deadly or benign, even though this should 52  
affect their initial diagnosis. In fact, sometimes payoffs will lead the decision-maker to choose 53  
the less probable alternative and this should be reflected by low confidence in the decision. 54

Little is known about how human observers adjust confidence in response to prior-payoff 55  
structures. In one perceptual study, the prior probabilities of target present versus absent 56  
affected the placement of the criteria for Type 1 and 2 judgments (Sherman et al., 2015), 57  
with some evidence that confidence better predicts performance for responses congruent with 58  
the more probable outcome than those that are incongruent. In the realm of social judg- 59  
ments, prior probabilities have been shown to modulate the degree of confidence, with higher 60  
confidence assigned to more probable outcomes (Manis et al., 1980). However, others have 61  
found counter-productive incorporation of priors, with over-confidence for low-probability 62  
outcomes and under-confidence for high-probability outcomes (Dunning et al., 1990). In 63  
regards to payoffs, early work on monetary incentives in perceptual categorization did col- 64  
lect confidence ratings, however they were not included in any analyses (Lee and Zentall, 65  
1966). Consideration of payoff structures is ubiquitous in animal studies of confidence that 66  
employ post-decisional wagering methods (Smith et al., 2003). For example, in the opt-out 67  
paradigm, to distinguish between low and high confidence, the animal chooses between a 68  
small, certain reward and a risky alternative with either high reward or no reward, for cor- 69  
rect and incorrect perceptual responses respectively (Kiani and Shadlen, 2009). However, 70  
because animals are motivated by their expected gain and not explicit verbal instructions, 71  
it is impossible to isolate decision confidence unconfounded with the subjective value of the 72  
reward. 73

Here, we seek to characterize how human observers adjust perceptual decisions and confi- 74  
dence in response to joint manipulation of priors and payoffs. First, we defined a normative 75

model of confidence judgments that factors in the prior-payoff structure of the environment. Then, we measured how well this model explains human behavior in an orientation-discrimination task, as compared to several sub-optimal decision models. We found that all observers made sub-optimal confidence judgments, but fell into two distinct groups depending on their strategy. These results highlight the importance of considering the effect of priors and payoffs on confidence, particularly in applied or real-world scenarios where they are likely to be non-uniform across the decision alternatives.

### 3 The Decision Models

In this section we describe the rationale and background for the modeling of Type 1 and Type 2 decision-making. We follow the example of a left-right orientation judgment followed by a binary low-high confidence judgment to match the experimental paradigm used in the present study. First the range of Type 1 models are identified, which assess the placement of the discrimination decision criterion under different prior-payoffs scenarios. Then the Type 2 models are outlined, describing the different potential relationships between the decision criteria for confidence and the criterion for discrimination.

#### 3.1 The Type 1 Decision

To make the Type 1 decision, observers must relate a noisy internal measurement,  $x$ , of the stimulus,  $s$ , where  $s \in \{s_L, s_R\}$ , to a binary response, which in the context of our experiment is “tilted left” (say “ $s = s_L$ ”) or “tilted right” (say “ $s = s_R$ ”). This is done by a comparison to an internal criterion,  $k_1$ , such that if  $x < k_1$ , the observer will respond with “tilted left”, and otherwise “tilted right” (Figure 1a). The only component of the Type 1 model where the observer has any control is deciding where to place the criterion. The optimal value of  $k_1$  ( $k_{opt}$ ) maximizes the expected gain, ensuring the observer makes the most points/money/etc. over the course of the experiment. The value of  $k_{opt}$  depends on three things:

- (i) The sensitivity of the observer,  $d'$ . In the standard model of the decision space,

$$P(x|s_L) \sim N(\mu_L, \sigma_L) \text{ and } P(x|s_R) \sim N(\mu_R, \sigma_R), \text{ with } \mu_L = -\mu_R \text{ and } \sigma_L = \sigma_R = 1.$$

Under this transformation, the sensitivity  $d'$  corresponds to the distance between the peaks of the two internal measurement distributions.

- (ii) The prior probability of each stimulus alternative,  $P(s_L)$  and  $P(s_R) = 1 - P(s_L)$ .
- (iii) The rewards for the four possible stimulus-response pairs,  $V_{r,s}$ , which are the rewards (positive) or costs (negative) of responding  $r$  when the stimulus is  $s$ .

An ideal observer that maximizes expected gain (Green and Swets, 1966) uses criterion

$$k_{opt} = \frac{\ln \beta_{opt}}{d'}, \quad (1)$$

where the likelihood ratio  $\beta_{opt}$  at the optimal criterion is a function of priors and payoffs:

$$\beta_{opt} = \frac{P(s_L) V_{L,L} - V_{L,R}}{P(s_R) V_{R,R} - V_{R,L}}. \quad (2)$$

In our experiment, 0 points are awarded for incorrect answers, allowing us to simplify:

$$\ln \beta_{opt} = \ln \frac{P(s_L) V_{L,L}}{P(s_R) V_{R,R}} = \ln \frac{P(s_L)}{P(s_R)} + \ln \frac{V_{L,L}}{V_{R,R}}. \quad (3)$$

Thus,  $k_{opt} = k_p + k_v$ , where  $k_p$  is the optimal criterion location if only priors were asymmetric and  $k_v$  is the optimal criterion if only the payoffs were varied. As can be seen in Eq. 3, the effects of priors and payoffs sum when determining the optimal criterion (illustrated in Figure 1b). When the priors are more similar, or the payoffs are closer to equal,  $k_{opt}$  is closer to the neutral criterion  $k_{neu} = 0$ . Note that in the case of symmetric payoffs,  $k_{opt}$  maximizes both expected gain and expected accuracy, whereas when asymmetric payoffs are involved,  $k_{opt}$  maximizes expected gain only (i.e.,  $k_{opt} \neq k_p$ ). This is because to maximize expected gain, from time to time the observer is incentivized to choose the less probable outcome because it is more rewarded.

## 3.2 Conservatism

Often, human observers use a sub-optimal value of  $k_1$  when the prior probabilities or payoffs are not identical for each alternative. A common observation is that the criterion is not

adjusted far enough from the neutral criterion towards the optimal criterion,  $k_{neu} < k_1 < k_{opt}$  122  
or  $k_{neu} > k_1 > k_{opt}$ , a behavior referred to as conservatism (Green and Swets, 1966; Maddox, 123  
2002). It is useful to express conservatism as a weighted sum of the neutral and optimal 124  
criterion: 125

$$k_1 = (1 - \alpha)k_{neu} + \alpha k_{opt} = \alpha k_{opt}, \quad (4)$$

with  $0 < \alpha < 1$  indicating conservative criterion placement. The degree of conservatism is 126  
greater the closer  $\alpha$  is to 0 (Figure 1c). Several studies have contrasted the conservatism for 127  
unequal priors versus unequal payoffs, typically finding greater conservatism for unequal pay- 128  
offs (Lee and Zentall, 1966; Ulehla, 1966; Healy and Kubovy, 1981; Ackermann and Landy, 129  
2015) with few exceptions (Healy and Kubovy, 1978). This may result from an underlying 130  
criterion-adjustment strategy that depends on the shape of the expected gain curve (as a 131  
function of criterion placement) and not just on the position of the optimal criterion max- 132  
imizing expected gain (Busemeyer and Myung, 1992; Ackermann and Landy, 2015) or a 133  
strategy that trades off between maximizing expected gain and maximizing expected accu- 134  
racy (Maddox, 2002; Maddox and Bohil, 2003). Given that the effects of priors and payoffs 135  
sum in Eq. 3, we will consider a sub-optimal model of criterion placement that has separate 136  
conservatism factors for payoffs and priors: 137

$$k_1 = \frac{1}{d'} \left[ \alpha_p \ln \frac{P(s_L)}{P(s_R)} + \alpha_v \ln \frac{V_{L,L}}{V_{R,R}} \right] = \alpha_p k_p + \alpha_v k_v. \quad (5)$$

The conservatism factors,  $\alpha_p$  and  $\alpha_v$ , scale these individually before they are summed to 138  
give the final conservative criterion placement, taking into account both prior and payoff 139  
asymmetries. This formulation allows for differing degrees of conservatism for priors and 140  
payoffs. 141

### 3.3 Type 1 Decision Models 142

We consider four models of the Type 1 discrimination decision in this paper, including the 143  
optimal model (i) and three sub-optimal models that include varying forms of conservatism 144  
(ii-iv): 145

$$(i) \Omega_{1,opt} : k_1 = k_{opt} = k_p + k_v \quad 146$$

$$(ii) \Omega_{1,1\alpha} : k_1 = \alpha k_{opt} = \alpha (k_p + k_v) \quad 147$$

$$(iii) \Omega_{1,2\alpha} : k_1 = \alpha_p k_p + \alpha_v k_v \quad 148$$

$$(iv) \Omega_{1,3\alpha} : \begin{cases} k_1 = \alpha_{pv} k_{opt} & \text{if } k_p \neq 0 \text{ and } k_v \neq 0 \text{ (i.e., both asymmetric)} \\ k_1 = \alpha_p k_p & \text{if } k_v = 0 \text{ (i.e., payoffs symmetric)} \\ k_1 = \alpha_v k_v & \text{if } k_p = 0 \text{ (i.e., priors symmetric)}. \end{cases} \quad 149$$

Thus, we consider models with no conservatism ( $\Omega_{1,opt}$ ), with an identical degree of conservatism due to asymmetric priors and payoffs ( $\Omega_{1,1\alpha}$ ), or different amounts of conservatism for prior versus payoff manipulations ( $\Omega_{1,2\alpha}$ ). In the fourth model, we drop the assumption (that was based on the optimal model) that effects of payoffs and priors on criterion sum, i.e., that behavior with asymmetric priors and payoffs can be predicted from behavior with each effect alone ( $\Omega_{1,3\alpha}$ ). We consider this final model because the additivity of criterion shifts (Eq. 3) has not yet been experimentally confirmed with human observers (Stevenson et al., 1990).

In all models, we also consider an additive bias term,  $\gamma$ , corresponding to a perceptual bias in perceived vertical. The bias is also included in the neutral criterion  $k_{neu} = \gamma$ . For clarity, however, we have omitted it from the mathematical descriptions of the models. Note that any observer best fit by  $\Omega_{1,opt}$  but with a  $\gamma$  significantly different from 0 would no longer be considered as having optimal behavior.

### 3.4 Confidence Criteria 163

Confidence judgments should reflect the belief that the selected alternative in the discrimination decision correctly matches the true world state. Generally speaking, the further the internal measurement is from a well-placed decision boundary, the more evidence there is for the discrimination judgment. This is instantiated in the extended SDT framework by the addition of two or more confidence criteria,  $k_2$  (Maniscalco and Lau, 2012, 2014). There are two such criteria for a binary confidence task and more confidence criteria when more than

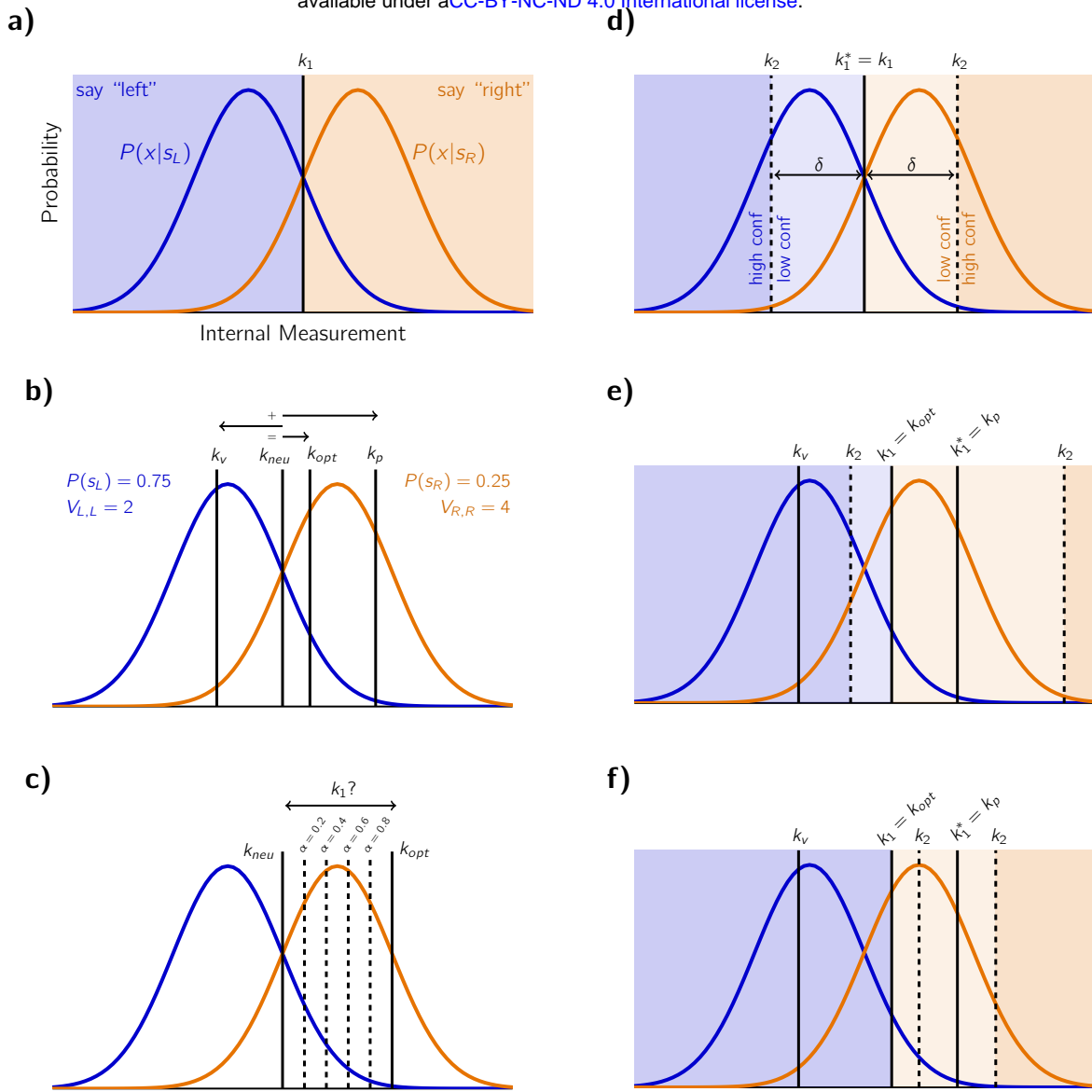


Figure 1: Illustration of the full SDT model. a) On each trial, an internal measurement of stimulus orientation is drawn from a Gaussian probability distribution conditional on the true stimulus value. The Type 1 criterion,  $k_1$ , defines a cut-off for reporting “left” or “right”. The ideal observer in a symmetrical priors and payoffs scenario is shown. b) The ideal observer’s criterion placement with both prior and payoff asymmetry. This prior asymmetry encourages a rightward criterion shift to  $k_p$  and the payoff asymmetry a leftward shift to  $k_v$ . The optimal criterion placement that maximizes expected gain,  $k_{opt}$ , is a sum of these two criterion shifts. For comparison, the neutral criterion,  $k_{neu}$  is shown. As the prior asymmetry is greater than the payoff asymmetry, 3:1 vs 1:2,  $k_{opt} \neq k_{neu}$ . c) A sub-optimal conservative observer will not adjust their Type 1 criterion far enough from  $k_{neu}$  to be optimal. The parameter  $\alpha$  describes the degree of conservatism, with values closer to 0 being more conservative and closer to 1 less conservative. d) In the case of symmetric payoffs and priors, the Type 2 confidence criteria,  $k_2$ , are placed equidistant from the Type 1 decision boundary, carving up the internal measurement space into a low- and high-confidence region for each discrimination response option. e) For the normative Type 2 model, the confidence criteria are placed symmetrically around a hypothetical Type 1 criterion that only maximizes accuracy ( $k_1^* = k_p$ ). This figure shows the division of the measurement space as per the prior-payoff scenario in (b). As a left-tilted stimulus is much more likely, this results in many high-confidence left-tilt judgments and few high-confidence right-tilt judgments. Note that left versus right judgments still depend on  $k_1$ . f) The same as in (e) but with small  $\delta$  value. Note the low-confidence region where confidence should be high (left of the left-hand  $k_2$ ). This happens because in this region the observer will choose the Type 1 response that conflicts with the accuracy-maximizing criterion, hence they will report low confidence in their decision. Note that the displacements of the criteria from the neutral criterion in this figure are exaggerated for illustrative purposes.



two confidence levels are provided. We restrict our treatment to the binary case, which can  
be trivially extended to include more gradations of confidence.

As illustrated in Figure 1d for the case of symmetric payoffs and priors, there is a  $k_2$   
confidence criterion on each side of the  $k_1$  decision boundary. If the measurement obtained  
is beyond one of these criteria relative to  $k_1$ , then the observer will report high confidence,  
and otherwise will report low confidence. Stated another way, the addition of the confi-  
dence criteria effectively divides the measurement axis into four regions: high-confidence  
left, low-confidence left, low-confidence right, and high-confidence right. The closer to the  
discrimination decision boundary that the observer places  $k_2$ , the more high-confidence re-  
sponses they will give. We denote this distance as  $\delta$ .  $\delta$  is not always assumed to be identical  
for both confidence criteria (e.g. Maniscalco and Lau, 2012), but we assumed a single value  
of  $\delta$  for model simplicity. Type 2 judgments were not incentivized in our experiment. Thus,  
there is no explicit cost function to constrain the distance parameter  $\delta$ , so the precise set-  
ting of  $\delta$  will not factor into the evaluation of how well the normative model fits observer  
behavior.

### 3.5 The Counterfactual Type 1 Criterion

The above description of how confidence responses are generated is well suited to cases where  
the payoffs are symmetric. This is because the optimal decision criterion maximizes both  
gain and accuracy. For an internal measurement at the discrimination boundary, it is equally  
probable that the stimulus had a rightward versus leftward orientation. Expressed another  
way, the log-posterior ratio at  $k_{opt}$  is 1. Thus, the distance from the discrimination boundary  
is a good measure for the probability that the Type 1 response is correct (i.e., confidence  
as we defined it above). This, however, is not the case when payoffs are asymmetric ( $k_1 =$   
 $k_p + k_v = k_{opt}$  where  $k_v \neq 0$ ), as the ideal observer maximizes gain but not accuracy. The  
log-posterior ratio is not 1 at  $k_{opt}$  but rather it is equal to 1 at  $k_p$ .

To extend the SDT model of confidence to asymmetric payoffs, we introduce a new  
criterion. The counterfactual criterion,  $k_1^*$ , is the criterion the ideal observer would have used  
if they ignored the payoff structure of the environment and exclusively maximized accuracy  
and not gain (i.e.,  $k_1^* = k_p$ ). It is around this criterion that the observer symmetrically places

confidence criteria in our normative model (Figure 1e). Whenever payoffs are symmetrical, 199  
 $k_1 = k_1^*$ . Figure 1f illustrates a situation unique to this model that may occur when payoffs 200  
are asymmetrical. Here, the value of  $\delta$  is sufficiently small that both  $k_2$  criteria fall on the 201  
same side of  $k_1$ . As a result, the region between  $k_1$  and the left-hand  $k_2$  criterion results in a 202  
low-confidence response despite being beyond the  $k_2$  boundary (relative to  $k_1^*$ ). This occurs 203  
because this region is to the right of  $k_1$  and thus, due to asymmetric payoffs, the observer will 204  
make the *less* probable choice, thus resulting in low confidence in that choice. Effectively, 205  
the left-hand confidence criterion is shifted from  $k_2$  to  $k_1$ . Here, we rely on the assumption 206  
that the confidence system is aware of the Type 1 decision (for further discussion of this 207  
issue, see Fleming and Daw, 2017). 208

The notion of an observer computing additional criteria for counterfactual reasoning is 209  
not new. For example, in the model of Type 1 conservatism of Maddox and Bohil (1998), 210  
where observers trade off gain versus accuracy,  $k_1$  is a weighted average of the optimal 211  
criteria for maximizing expected gain ( $k_{opt}$ ) and for exclusively maximizing accuracy ( $k_p$ ). 212  
In Zylberberg et al. (2018), observers learned prior probabilities of each stimulus type by an 213  
updating decision-making mechanism that computes the confidence the observer would have 214  
had if they had used the neutral criterion ( $k_{neu}$ ) for their Type 1 judgment. We suggest that 215  
for determining confidence in the face of asymmetric payoffs, optimal observers compute 216  
the confidence they would have reported if they had instead used the  $k_p$  criterion for the 217  
discrimination judgment. 218

### 3.6 Type 2 Decision Models 219

In addition to the normative model we just described (i), we considered four sub-optimal 220  
models (ii-v) for the counterfactual Type 1 criterion about which the Type 2 criteria are 221  
symmetrically arranged: 222

(i)  $\Omega_{2,acc} : k_1^* = k_p$  223

(ii)  $\Omega_{2,acc+cons} : k_1^* = \alpha_p k_p$  224

(iii)  $\Omega_{2,gain} : k_1^* = k_{opt}$  225

(iv)  $\Omega_{2,gain+cons} : k_1^* = k_1$  226

(v)  $\Omega_{2,neu} : k_1^* = k_{neu}$  227

All of these models are characterized by the placement of the counterfactual criterion,  $k_1^*$ ; 228  
the distance  $\delta$  is the only free parameter for all models. In the normative model ( $\Omega_{2,acc}$ ), the 229  
confidence criteria are systematically shifted so that they are centered on the discrimination 230  
criterion that maximizes accuracy. We also consider a model in which confidence criteria 231  
are centered on the criterion that maximizes gain ( $\Omega_{2,gain}$ ), which is incorrect behavior in 232  
the case of asymmetric payoffs. In the neutral model ( $\Omega_{2,neu}$ ), confidence criteria remain 233  
fixed around the neutral Type 1 criterion regardless of the prior or payoff manipulation. 234  
Finally, for the models that shift in response to priors and payoffs, we consider that con- 235  
servatism in the discrimination criterion placement also affects  $k_1^*$ , either for the accuracy 236  
model ( $\Omega_{2,acc+cons}$ ) or the gain model ( $\Omega_{2,gain+cons}$ ). In the latter model,  $k_1^*$  is identical to 237  
 $k_1$ . For the other models, some combinations of priors and payoffs will decouple  $k_1^*$  from  $k_1$ . 238  
For the  $\Omega_{2,acc+cons}$  model, the decoupling only occurs for asymmetric payoffs. For the other 239  
models, this decoupling occurs whenever priors or payoffs are asymmetric. 240

Our models assume that  $\delta$  are placed symmetrically around  $k_1^*$ . However, the ability to 241  
identify the underlying Type 2 model will not be affected by this assumption. Consider an 242  
observer whose low-confidence region to the left of  $k_1^*$  was always greater than their low- 243  
confidence region to the right of  $k_1^*$ , such that  $k_1^* - k_{2-} > k_{2+} - k_1^*$ . Then, the estimate of  $\delta$  244  
would be similar because the experiment design tested the mirror prior-payoff condition (i.e., 245  
for fixed  $k_2$ , one condition would have  $k_1^*$  attracted to neutral and the other repelled, which is 246  
not the behaviour of  $k_1^*$  in any Type 2 model). Thus, the best-fitting model would be unlikely 247  
to change when  $\delta$  is asymmetric, but it would fit less well. Alternatively, an asymmetry in 248  
 $\delta$  could be mirrored about the neutral criterion (e.g., the low confidence region closest to 249  
the neutral criterion is always smaller). Then, the  $\delta$  asymmetry would be indistinguishable 250  
from a bias in the conservatism parameters. Ultimately, the confidence criteria are yoked to 251  
 $k_1^*$ , and it is the patterns of criteria shift from all conditions jointly that are captured by the 252  
model comparison. 253

## 4 Methods 254

### 4.1 Participants 255

Ten participants (5 female, age range 22-43 years, mean 27.0 years) took part in the experiment. All participants had normal or corrected-to-normal vision, except one amblyopic participant. All participants were naive to the research question, except for three of the authors who participated. On completion of the study, participants received a cash bonus in the range of \$0 to \$20 based on performance. In accordance with the ethics requirements of the Institutional Review Board at New York University, participants received details of the experimental procedures and gave informed consent prior to the experiment. 256  
257  
258  
259  
260  
261  
262

### 4.2 Apparatus 263

Stimuli were presented on a gamma-corrected CRT monitor (Sony G400, 36 x 27 cm) with a 1280 x 1024 pixel resolution and an 85 Hz refresh rate. The experiment was conducted in a dimly lit room, using custom-written code in MATLAB version R2014b (The MathWorks, Natick, MA), with PsychToolbox version 3.0.11 (Brainard, 1997; Pelli, 1997; Kleiner et al., 2007). A chin-rest was used to stabilize the participant at a viewing distance of 57 cm. Responses were recorded on a standard computer keyboard. 264  
265  
266  
267  
268  
269

### 4.3 Stimuli 270

Stimuli were Gabor patches, either right (clockwise) or left (counterclockwise) of vertical, presented on a mid-gray background at the center of the screen. The sinusoidal grating had a spatial frequency of 2 cycle/deg, a peak contrast of 10%, and a Gaussian envelope (SD: 0.5 deg). The phase of the grating was randomized on each trial to minimize contrast adaptation. 271  
272  
273  
274  
275

### 4.4 Experimental Design 276

Orientation discrimination (Type 1, 2AFC, left/right) and confidence judgments (Type 2, 2AFC, low/high) were collected for seven conditions defined by the prior and payoff struc- 277  
278

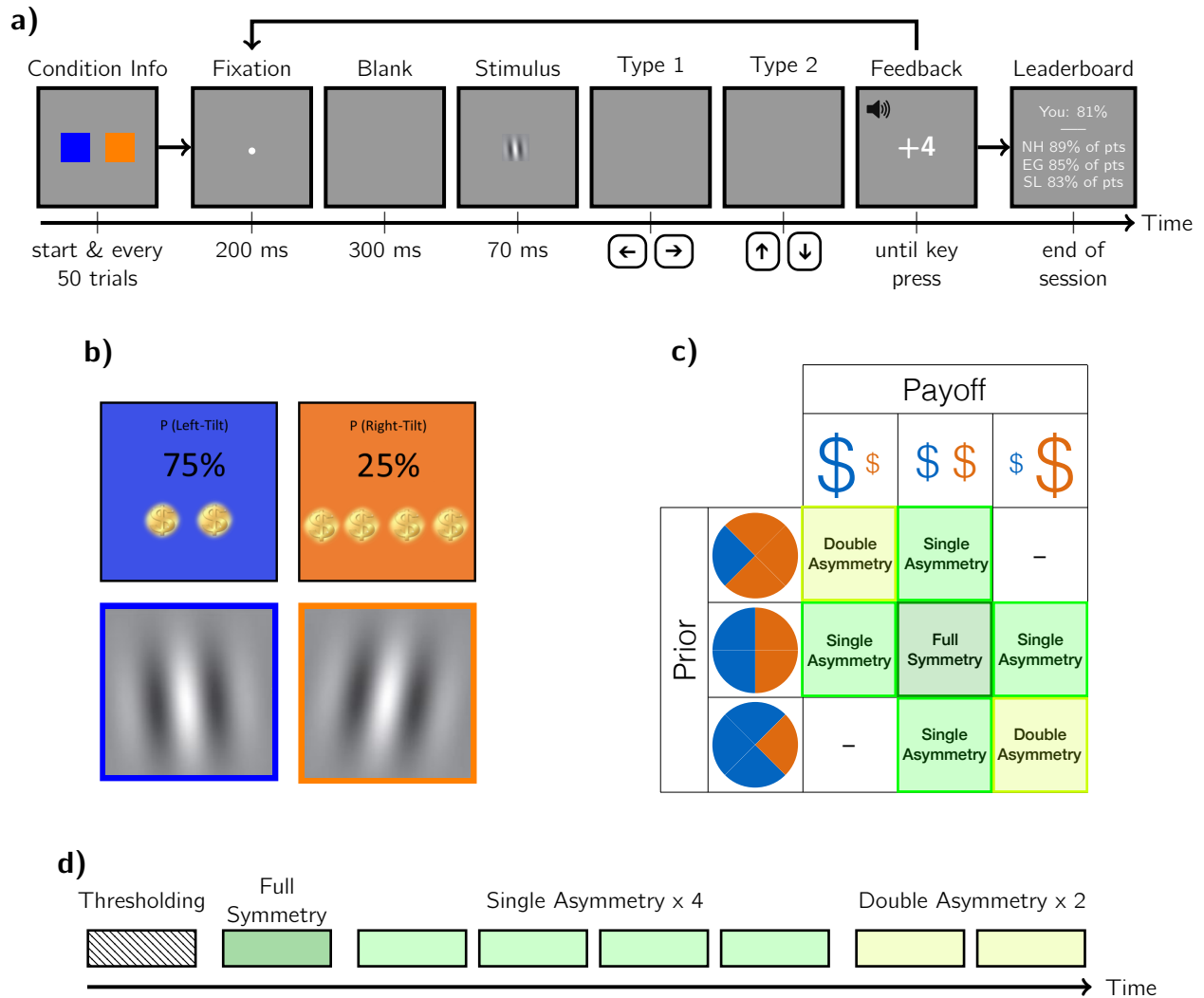


Figure 2: Experimental methods. a) Trial sequence including an outline of the initial condition information screen (see part (b) for details) and final (mock) leaderboard screen. Participants were shown either a right- or left-tilted Gabor and made subsequent Type 1 and Type 2 decisions before being awarded points and given auditory feedback based on the Type 1 discrimination judgment. b) Sample condition information displays from a double-asymmetry condition. Below: Example Gabor stimuli, color-coded blue for left- and orange for right-tilted. The exact stimulus orientations depended on the the participant's sensitivity. c) Condition matrix. Pie charts show the probability of stimulus alternatives (25, 50, or 75%) and dollar symbols represent the payoffs for each alternative (2, 3, or 4 pts). Squares are colored and labeled by the type of symmetry. d) Timeline of the eight sessions. The order of conditions was randomized within the single- and within the double-asymmetry conditions.

ture. The probability of a right-tilted Gabor could be 25, 50, or 75%. The points awarded 279  
for correctly identifying a right- versus a left-tilt could be 4:2, 3:3, or 2:4. In the 3:3 pay- 280  
off scheme, a correct response was awarded 3 points. In the 2:4 and 4:2 schemes, correct 281  
responses were awarded 2 or 4 points depending on the stimulus orientation. Incorrect re- 282  
sponses were not rewarded (0 points). The prior and payoff structure was explicitly conveyed 283  
to the participant before the session began (Fig. 2b) and after every 50 trials. Condition 284  
order was randomized within condition type (Fig. 2c): no asymmetry (50%, 3:3), single 285  
asymmetry (50%, 4:2; 50%, 2:4; 25%, 3:3; 75%, 3:3), or double asymmetry (25%, 4:2; 75%, 286  
2:4). Note that two of the possible double asymmetry conditions (25%, 2:4; and 75%, 4:2) 287  
were not tested because these conditions incentivized one response alternative to such a de- 288  
gree that they would not be informative for model comparison. Participants first completed 289  
the full-symmetry condition, followed by the single-asymmetry conditions in random order, 290  
and finally the double-asymmetry conditions, also in random order (Fig. 2d). Each condition 291  
was tested in a separate session with no more than one session per day. 292

## 4.5 Thresholding Procedure 293

A thresholding procedure was performed prior to the main experiment to equate difficulty 294  
across observers to approximately  $d' = 1$ . Observers performed a similar orientation discrim- 295  
ination judgment as in the main experiment. Absolute tilt magnitude varied in a series of 296  
interleaved 1-up-2-down staircases to converge on 71% correct. Each block consisted of three 297  
staircases with 60 trials each. Participants performed multiple blocks until it was determined 298  
that performance had plateaued (i.e., learning had stopped). Preliminary thresholds were 299  
calculated using the last 10 trials of each staircase. At the end of each block, if none of the 300  
three preliminary thresholds were better than the best of the previous block's preliminary 301  
thresholds, then the stopping rule was met. As a result, participants completed a minimum 302  
of two blocks and no participant completed more than five blocks. A cumulative Gaussian 303  
psychometric function was fit by maximum likelihood to all trials from the final two blocks 304  
(360 trials total). The slope parameter was used to calculate the orientation corresponding 305  
to 69% correct for an unbiased observer ( $d' = 1$ ; Macmillan and Creelman, 2005). This 306  
orientation was then used for this subject in the main experiment. Thresholds ranged from 307

0.36 to 0.78 deg, with a mean of 0.59 deg.

308

## 4.6 Main Experiment

309

Participants completed seven sessions, each of which had 700 trials with the first 100 treated as warm-up and discarded from the analysis. All subjects were instructed to hone their response strategy in the first 50 trials to encourage stable criterion placement. The trial sequence is outlined in Fig. 2a. Each trial began with the presentation of a fixation dot for 200 ms. After a 300 ms inter-stimulus interval, a Gabor stimulus was displayed for 70 ms. Participants judged the orientation (left/right) and then indicated their confidence in that orientation judgment (high/low). Feedback on the orientation judgment was provided at the end of the trial by both an auditory tone and the awarding of points based on the session's payoff structure. Additionally, the running percentage of potential points earned was shown on a leaderboard at the end of each session to foster inter-subject competition. Participants' cash bonus was calculated by selecting one trial selected at random from each session and awarding the winnings from that trial, with a conversion of 1 point to \$1, capped at \$20 over the sessions. Total testing time per subject was approximately 8 hrs.

310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322

## 4.7 Model Fitting

323

Detailed description of the model-fitting procedure can be found in Supplementary Information (Sections 1 and 2). Briefly, model fitting was performed in three sequential steps: fitting of  $d'$ , Type 1 models, then Type 2 decisions. Each session provided one measurement of  $d'$ , which we used in a hierarchical Bayesian model to estimate each participant's true underlying  $d'$  value across the entire experiment. For Type 1 and Type 2 models, we calculated the log likelihood of the data given a dense grid of parameters (e.g.,  $\alpha$ ,  $\gamma$ , and  $\delta$ ) using multinomial distributions defined by the stimulus type, discrimination response, and confidence response. All seven conditions were fit jointly. We then calculated model evidence by marginalizing over all parameter dimensions and then normalizing to account for grid spacing.

324  
325  
326  
327  
328  
329  
330  
331  
332  
333

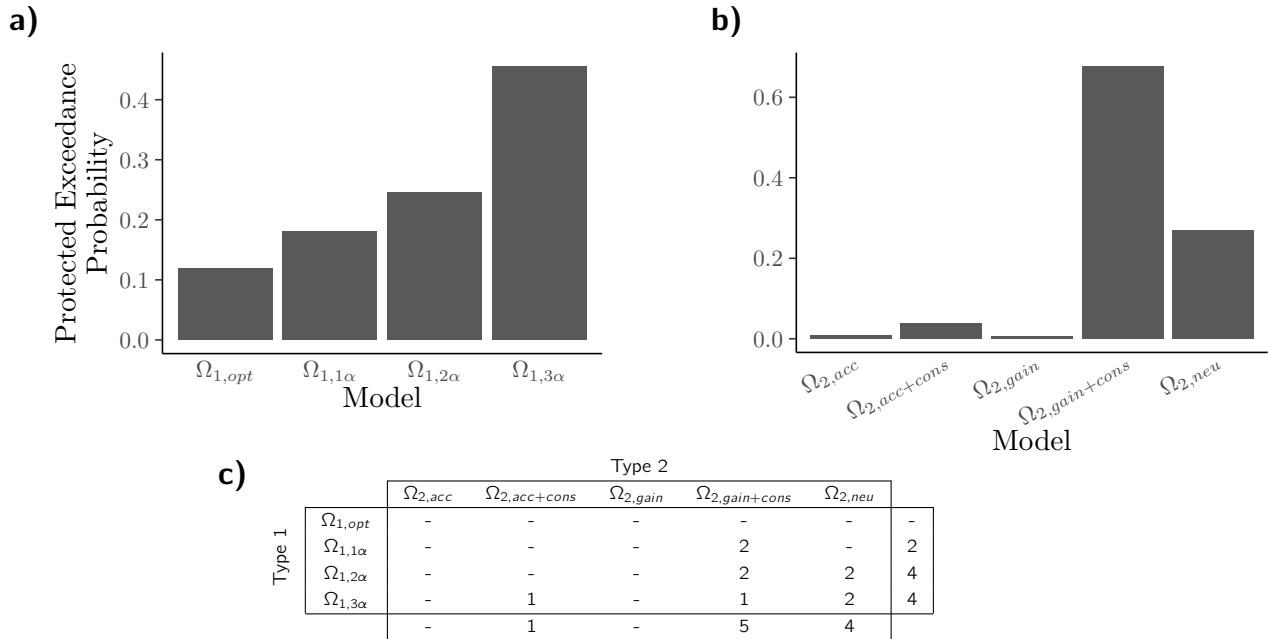


Figure 3: Model comparison for the Type 1 and Type 2 responses. a) The protected exceedance probabilities (PEPs; see text for details) of the four Type 1 models. b) PEPs of the five Type 2 models. Note that model comparisons were performed first for Type 1 and then for Type 2 responses, using the best-fitting Type 1 model and parameters, on a per-subject basis, in the Type 2 model evaluation. c) Best-fitting models for each participant.

## 5 Results

334

We sought to understand how observers make perceptual decisions and confidence judgments in the face of asymmetric priors and payoffs. Participants performed an orientation-discrimination task followed by a confidence judgment. To account for the behavior, we defined two sets of models, which were fit in a two-step process. Type 1 models defined the contribution of conservatism to the discrimination responses. Type 2 models defined the role of priors and payoffs in the confidence reports.

335

336

337

338

339

340

### 5.1 Model Fits

341

Type 1 models were first fit using the discrimination responses alone. Four models were compared: optimal criterion placement ( $\Omega_{1,opt}$ ), equal conservatism for priors and payoffs ( $\Omega_{1,1\alpha}$ ), different degrees of conservatism for priors and payoffs ( $\Omega_{1,2\alpha}$ ), and a model in which there was a failure of summation of criterion shifts in the double-asymmetry condition ( $\Omega_{1,3\alpha}$ ). Fitting the Type 1 models also provided an estimate of response bias,  $\gamma$ . We performed a Bayesian model selection procedure using the SPM12 Toolbox (Wellcome Trust Centre for

342

343

344

345

346

347



Neuroimaging, London, UK) to calculate the protected exceedance probabilities (PEPs) 348  
for each model (Figure 3a). The exceedance probability (EP) is the probability that a 349  
particular model is more frequent in the general population than any of the other tested 350  
models. The PEP is a conservative measure of model frequency that takes into account the 351  
overall ability to reject the null hypothesis that all models are equally likely in the population 352  
(Stephan et al., 2009; Rigoux et al., 2014). Overall, an additional parameter in the double- 353  
asymmetry conditions was needed to explain Type 1 criterion placement, indicating a failure 354  
of summation of criterion shifts (i.e., the best-fitting model was  $\Omega_{1,3\alpha}$ ). 355

In the second step, the Type 2 models were fit using each participant's best Type 1 model 356  
and the associated maximum *a posteriori* (MAP) parameter estimates. The Type 2 models 357  
differed in the placement of the Type 2 criteria, which split the internal response axis into 358  
“high” and “low” confidence regions, for each “right” and “left” discrimination response. 359  
We modeled the two Type 2 criteria as shifting to account for only the prior probability, 360  
maximizing accuracy with the confidence response ( $\Omega_{2,acc}$ ; the normative model), shifting 361  
the confidence criteria in response to payoff manipulations ( $\Omega_{2,gain}$ ; a sub-optimal model), or 362  
failing to move the confidence criteria away from neutral at all ( $\Omega_{2,neu}$ ; a sub-optimal model). 363  
We also tested models where the conservatism found in the Type 1 decisions carried over into 364  
the confidence decision ( $\Omega_{2,acc+cons}$  and  $\Omega_{2,gain+cons}$ ; both sub-optimal). We again compared 365  
the models quantitatively with PEPs (Figure 3b). The favored model,  $\Omega_{2,gain+cons}$ , shifts 366  
the confidence criteria in response to both prior and payoff manipulations. Furthermore, 367  
the conservatism that participants exhibited in the Type 1 decisions carried over into the 368  
placement of the confidence criteria. 369

Figure 3c shows the best-fitting models for individual participants, according to the 370  
amount of relative model evidence (here the marginal log-likelihood). Each of the Type 1 371  
models except the optimal ( $\Omega_{1,opt}$ ) was a best-fitting model for at least one of the ten partici- 372  
pants. Similarly, no one was best fit by the normative Type 2 model either ( $\Omega_{2,acc}$ ). Overall, 373  
there is no clear pattern between the pairings of Type 1 and Type 2 models. 374

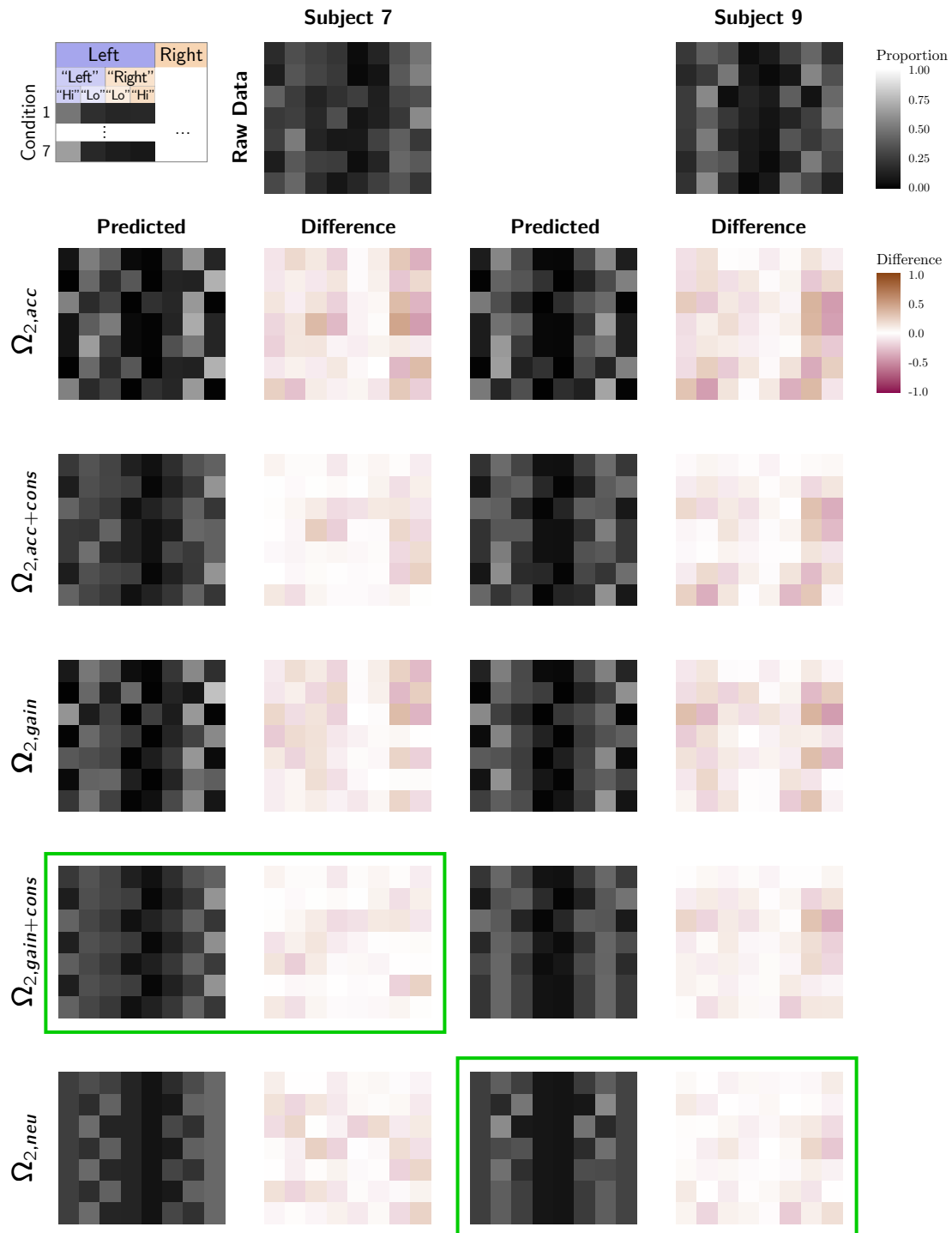


Figure 4: Visualization of the raw and predicted response rates for two example participants. Grids are formed of the seven conditions (rows) and the eight possible stimulus-response-confidence combinations (columns). See Figure S3 in the Supplement for condition order. The fill indicates the proportion of trials for that condition and stimulus that have that combination of response and confidence. Top row: Raw response rates of two example subjects. Subsequent rows, columns 1 and 3: Predicted response rates for each Type 2 model using the best-fitting parameters of the best-fitting Type 1 model for that individual. Columns 2 and 4: Difference between raw and predicted response rates. Green boxes: winning models (Subject 7:  $\Omega_{gain+cons}$ ; subject 9:  $\Omega_{neu}$ ).

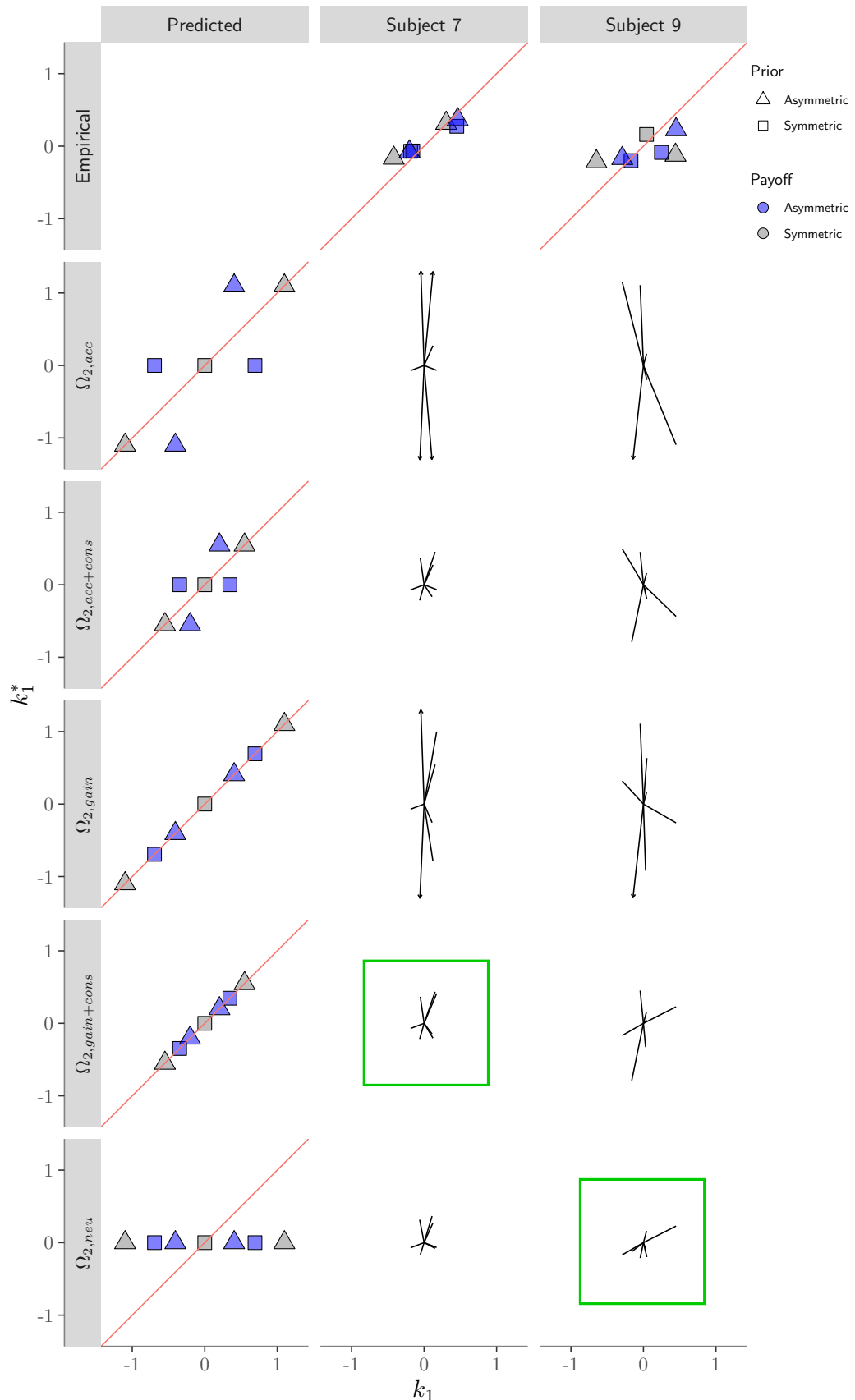


Figure 5: Comparison of the empirical and predicted  $k_1$  and  $k_1^*$ . Top row: empirical criteria of two example observers. The  $k_1^*$  was calculated as the midpoint between the two empirical  $k_2$  (see Figure S1 for  $k_2$  calculation details). Left column: predicted relationship between the Type 1 and Type 2 criteria ( $d' = 1$ ; either  $\Omega_{1,opt}$  or  $\Omega_{1,1\alpha}$  with  $\alpha = 0.5$ ). Grey and square symbols: symmetry conditions. Triangles: prior asymmetry. Blue symbols: payoff asymmetry. Polar plots: residuals between empirical data and model prediction based on best-fitting parameters, plotted as vectors. Arrowheads: residuals greater than plot bounds.

## 5.2 Model Checks

We performed several checks on the fitted data to ensure that parameters were capturing expected behavior and that the models could predict the data well (reported in detail in Section 3 of the Supplementary Information). The quality of a model is not only dependent on how much more likely it is than others, but it is also dependent on its overall predictive ability. To visualize each model’s ability to predict the proportion of each response type (“right” vs. “left” x “high” vs. “low”), we calculated the expected proportion of each response type given the MAP parameters for each model and participant. We compared the predicted response proportions to the empirical proportions (Figure 4). Larger residuals are represented by more saturated colors. For the best-fitting models, the residuals are small and unpatterned.

We also compared the Type 1 criteria and the counterfactual confidence criteria (Figure 5). We constrained the empirical counterfactual confidence criterion to be the midpoint between the two Type 2 criteria (i.e.,  $k_1^* \equiv (k_{2-} + k_{2+})/2$ ). Using  $k_1^*$ , the predictions made by the Type 2 models are highly distinguishable. In the left-most column, predicted  $k_1$  and  $k_1^*$  for each session are shown for each model, assuming  $d' = 1$  and either  $\Omega_{1,opt}$  or  $\Omega_{1,1\alpha}$  where  $\alpha = 0.5$ . In the top row, empirical criteria from the same two example participants as in Figure 4 are shown. Empirical criteria are calculated with the standard SDT method (detailed in Section 1 of the Supplementary Information, see Figure S1).

The visualization in the top row and left-most column of Figure 5 illustrates several behavioral phenomena. The response bias,  $\gamma$ , results in a shift in all criteria in the same direction, translating all data points parallel to the identity line. Conservatism is represented by an attraction of all data toward the origin on the x-axis for Type 1 and the y-axis for Type 2 judgments. The Type 2 models predict qualitatively different arrangements of the data points. If the prior and payoff asymmetries affect the placement of the Type 1 criterion but not the Type 2 criteria ( $\Omega_{2,neu}$ ), the data are clustered along a single value on the y-axis. If the prior and the payoff affect the placement of the Type 1 and Type 2 criteria equally ( $\Omega_{2,gain}$ ), then the data fall on the identity line. With normative behavior ( $\Omega_{2,acc}$ ), the prior asymmetry conditions (grey triangles) fall on the identity line because confidence tracks the prior, while in the payoff asymmetry conditions (blue squares), the data have the same  $k_2$

midpoint as in the neutral condition (grey squares) because confidence does not track the  
payoff.

Vectors in all 10 of the bottom right polar plots represent the difference (i.e., the residual)  
between the empirical and the predicted criteria from the model fits. While the model  
prediction column is based on fixed parameters, the predicted data used for the 10 polar  
plots uses parameters that best fit the participant's data using that model. It is immediately  
clear that the normative model (second row) does a poor job of describing participants'  
behavior, and that conservatism is a necessary component of the models.

### 5.3 Conservatism

We first measured the relative magnitude of conservatism due to priors and payoffs. Figure 6a  
shows fitted  $\alpha_p$  and  $\alpha_v$  under the most complex conservatism model ( $\Omega_{1,3\alpha}$ ) and Figure 6b  
shows them under the best-fitting model for each observer. In these figures, eight of the  
ten participants were conservative in their criterion placement for both prior and payoff  
manipulations, as indicated by data points in the gray regions. Of the eight participants  
that displayed conservatism, five were significantly more conservative for payoff asymmetries  
than prior asymmetries ( $\alpha_v < \alpha_p$ ), whereas only one was significant in the opposite direction  
( $\alpha_p < \alpha_v$ ). At the group level, however, we did not find a significant difference between the  
best fitting  $\alpha_v$  and  $\alpha_p$ , either for the best-fitting Type 1 model or the winning model (paired  
t-tests,  $p > 0.05$ ). Note that the negative  $\alpha$  values derive from a participant who shifted  
criteria consistently in the opposite direction expected from a rational observer in response  
to manipulations of payoffs and priors.

An additional implication of SDT is that an ideal observer's criterion shift due to pay-  
offs and due to priors should sum when both asymmetries are present as in Figure 1b:  
 $k_{pv} = k_p + k_v$  (Stevenson et al., 1990). Figure 6c shows the prediction of this additive rule.  
Although the difference between the predicted and actual criterion shift is marginally signif-  
icant ( $t = 2.41, p = .039$ ), this effect is driven by the four observers best fit by  $\Omega_{1,3\alpha}$ . Each of  
these four observers had 95% CIs that did not overlap with the identity line. We show the  
criterion placement in the double-asymmetry cases in Figure 6d. Most observers did not shift  
their criterion far enough from neutral to the optimal placement,  $k_{opt}$ . Three observers, how-

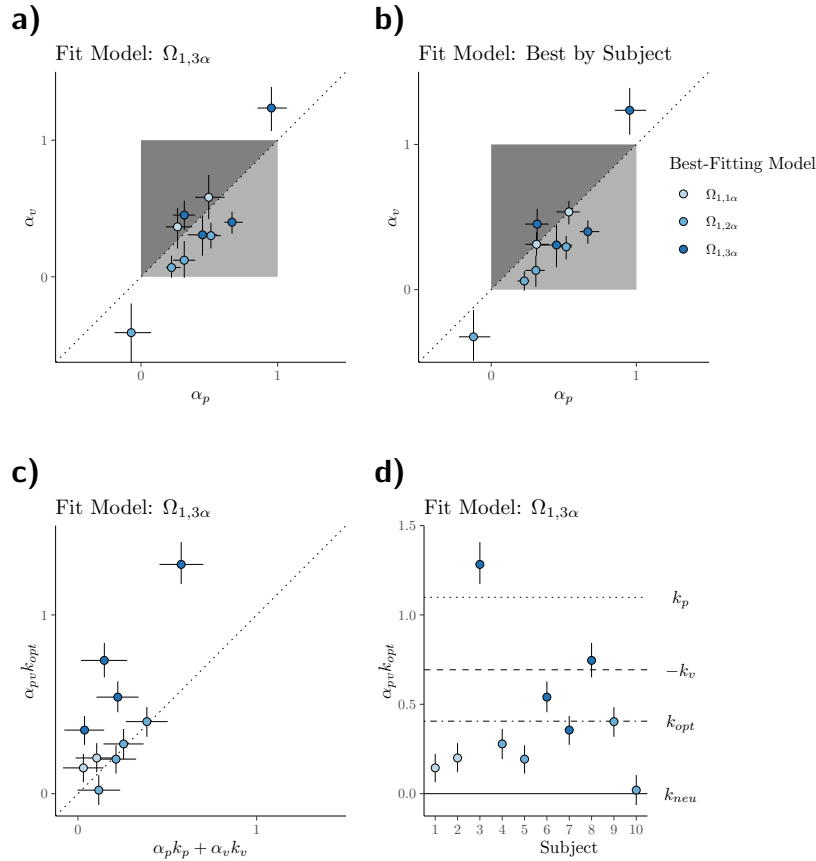


Figure 6: Conservatism for Type 1 decision making. a) A comparison of the extent of conservatism under payoff versus prior asymmetries. Each data point represents the best-fitting conservatism parameters of a single observer when fit by  $\Omega_{1,3\alpha}$ . These parameters are only contingent on the conservatism in the single-asymmetry conditions. In this model, conservatism in the double-asymmetry conditions is captured by a separate model parameter. Darker marker fill: additional conservatism parameters were required to fit to that observer's data. Dashed line: equality line. Dark grey region: conservatism greater for prior than payoff manipulations (i.e.,  $\alpha_p < \alpha_v$ ). Light grey region: conservatism is greater for payoffs (i.e.,  $\alpha_p > \alpha_v$ ). Data points outside these regions are not consistent with conservative criterion placement. b) Same as (a) using fit parameters from the best-fitting Type 1 model for each observer. c) Test of summation of criterion shifts using the  $\Omega_{1,3\alpha}$  model fits. Observers who required a third  $\alpha$  to capture their data (i.e., were best fit by  $\Omega_{1,3\alpha}$ ) had criterion shifts for the double-asymmetry conditions that were not well predicted as the sum of the shifts in the single-asymmetry conditions. d) Criterion placement in the double-asymmetry conditions. These are the same data as in the y-axis of (c), but extended to more easily compare the actual criterion placement with potential other task-relevant criteria. Horizontal criteria lines assume  $d'=1$ .

ever, placed their criterion beyond  $k_{opt}$ , with two stopping short of the accuracy-maximizing 434  
criterion  $k_p$ . 435

In summary, we find that conservatism for priors and conservatism for payoffs do not sum, 436  
as traditional SDT predicts. Conservatism applied to priors and payoffs in the discrimination 437  
decision was also incorporated into the confidence decision. Participants further deviated 438  
from normative behavior by shifting their confidence criteria in response to asymmetric 439  
payoffs, which do not inform the probability of a correct discrimination response. 440

## 6 Discussion 441

### 6.1 Type 1 Judgments 442

We conducted an orientation-discrimination task in which the prior probability of response 443  
alternatives and the payoff matrix varied across sessions. Binary confidence reports were 444  
collected after each discrimination judgment to gauge the observer's subjective appraisal of 445  
the probability they were correct in their judgment. Observers were found to be conserva- 446  
tive in the placement of the discrimination criterion,  $k_1$ , as revealed by the Type 1 model 447  
comparison. Instead of placing the criterion at the optimal location, as determined by the 448  
priors and payoffs, they had the tendency to place  $k_1$  between the optimal criterion and the 449  
neutral criterion. While we did find evidence of different degrees of conservatism for payoff 450  
versus prior asymmetries at the individual-subject level, we found no evidence at the group 451  
level that conservatism was stronger when the payoffs were asymmetrical than when the pri- 452  
ors were asymmetrical. Differences in conservatism were more apparent in previous studies 453  
(Lee and Zentall, 1966; Ulehla, 1966; Healy and Kubovy, 1981; Maddox, 2002; Ackermann 454  
and Landy, 2015), but not all (Healy and Kubovy, 1978). Several factors may be contribut- 455  
ing to the observed conservatism of individual observers. Candidate explanations include 456  
the hypothesis that observers trade off between maximizing gains and maximizing accuracy 457  
(Maddox and Bohil, 1998), as it may be hard for the observer to sacrifice accuracy for ex- 458  
pected gain. Alternatively, conservatism could depend on the criterion-adjustment strategy 459  
(Busemeyer and Myung, 1992), which may be differentially influenced by subjective factors 460  
such as subjective probability and subjective utility (Ackermann and Landy, 2015). This 461

explanation suggests that it is effortful to shift the criterion far from the neutral criterion 462  
for an inconsequential gain. Another possibility is that it may be a combination of the two, 463  
as suggested by Maddox and Bohil (2003). 464

In additional analyses we explored the nature of conservatism further, both by fitting 465  
Type 1 models with varying levels of complexity as well as testing the predictions of several 466  
possible models for the case of both prior and payoff asymmetries. All participants were best 467  
described by a model with some form of conservatism, with the majority best fit with two 468  
or three separate conservatism parameters. In the extreme case, where three conservatism 469  
parameters were needed, we find cases where additivity of criterion shifts was not obtained, 470  
as predicted by Healy and Kubovy (1981). By additivity we mean that the criterion shifts 471  
induced by priors or payoffs sum when both are present. In our sample population, additivity 472  
was not found for 40% of observers, which provides a similarly inconclusive follow-up to 473  
previous attempts at testing additivity (Stevenson et al., 1990). Yet, the Bayesian Model 474  
Selection procedure indicated that this was the winning model. Taking into account the 475  
evidence for each model, as well as penalizing model complexity, we find that the most 476  
complex Type 1 model does the best job of describing the behavior of our sample population. 477  
Without any sizeable, systematic deviation from additivity (Figure 6c), it is reasonable to 478  
suggest this third conservatism parameter is capturing something else, relating to strategy 479  
or noise, on the part of the observers. 480

How consistent is additivity with the various explanations of unequal conservatism for 481  
priors and payoffs? In Sect. 1.4 of the Supplementary Information, we demonstrate that the 482  
gain-accuracy trade-off strategy is equivalent to our  $\Omega_{1,2\alpha}$  model, for which additivity holds. 483  
Therefore, observers best fit by this model may be simply trading off between maximizing 484  
gain and maximizing accuracy. Turning to the criterion-adjustment strategy explanation of 485  
conservatism, behavior might deviate from additivity depending on whether conservatism, 486  
which acts as a scale factor on criterion shifts, is applied before or after the individual shifts 487  
for priors and the payoffs are combined. If conservatism is applied to these components 488  
individually, and then the resulting criteria are summed, this is equivalent to the  $\Omega_{1,1\alpha}$  or 489  
 $\Omega_{1,2\alpha}$  models. If, however, the criterion is adjusted after both the priors or payoffs have been 490  
applied, then the rate of change in reward based on the objective or subjective gain functions 491



is in no way constrained to match that of the single-asymmetry cases. Yet, we found that the discrimination criterion in the double-asymmetry cases was placed beyond the criterion for 30% of observers, which is not consistent with a reluctance to shift the criterion sufficiently from neutral. In fact, these criteria are biased in the direction of the accuracy-maximizing criterion, as would be expected under the gain-accuracy trade-off hypothesis. However, we cannot distinguish the trade-off hypothesis from a liberal criterion placement in the double-asymmetry case because, in our task, the prior odds ratio was always more asymmetric than the rewards ratio, always placing the optimal criterion on the same side of the neutral criterion as the accuracy-optimizing criterion.

So far, we have considered explanations of conservatism that are a result of prior and payoff factors. An alternate metacognitive source of conservatism proposed by Kubovy (1977) implicates the  $d'$  component of Eq. 5. Observers likely form an estimate of their overall performance from experience with the task. If they happen to overestimate performance (i.e.,  $\hat{d}' > d'$ ), then it follows from Eq. 5 that  $k_1 < k_{opt}$ , and vice versa for underestimation. Note that this is not a form of confidence in the response for a given trial, but a more general metacognitive appraisal of the difficulty of the task. According to this hypothesis, most of the observers would have been overestimating performance, with the one observer with liberal criterion placement underestimating their performance. While it is not uncommon to find overestimation of performance in the metacognitive literature (Mamassian, 2008, 2016), this explanation alone is insufficient as we find differences in the degree of conservatism for priors versus payoffs for some participants. Thus, we conclude that the conservatism observed in this task is likely due to a combination of possible factors, including noisy behavior, strategies to trade off gain versus accuracy, sub-optimal criterion adjustment, and bias in participants' judgments of their own  $d'$ .

## 6.2 Type 2 Judgments

We now turn to the Type 2 results, i.e., how observers form confidence judgments about the discrimination decision. Five Type 2 models were characterized by the placement rule for the counterfactual Type 1 criterion,  $k_1^*$ , around which the confidence criteria,  $k_2$ , were symmetrically placed. We tested whether this counterfactual criterion coincided with the

accuracy-maximizing criterion, the gain-maximizing criterion, either of these options with 521  
the Type 1 conservatism applied, or whether it remained fixed at the neutral criterion. We 522  
found no observer was best fit by the accuracy-maximizing optimal model ( $\Omega_{2,acc}$ ), with 523  
the majority split between the gain-maximizing-with-conservatism model ( $\Omega_{2,gain+cons}$ ) or 524  
the fixed neutral model ( $\Omega_{2,neu}$ ). One participant was best fit by the accuracy-maximizing 525  
with conservatism model ( $\Omega_{2,acc+cons}$ ). Overall, Bayesian model selection favored the gain- 526  
maximizing model with conservatism. When considering the best-fitting Type 1 model, we 527  
find no clear pattern between the number of conservatism parameters required to explain 528  
behavior and the placement strategy for confidence criteria. 529

We first turn our focus to the subset of observers who were best fit by the model in which 530  
confidence criteria remained fixed around neutral ( $\Omega_{2,neu}$ ). In this model, the perceived 531  
magnitude of the tilt was all that was used to compute confidence. These observers correctly 532  
did not allow the payoff structure of the environment to affect confidence, unlike the other 533  
top-winning model. However, it is sub-optimal not to include the additional information 534  
provided via the priors for the response alternatives but the lack of adaptability should not 535  
be taken as evidence of an inability to adapt. It is possible that these observers ignored the 536  
prior-payoff structure entirely for confidence, and instead opted for a criterion-placement 537  
strategy that would work best for all conditions of the experiment. Future experiments 538  
could incentivize accurate confidence judgments to test this hypothesis. 539

In the winning Type 2 model, observers placed  $k_1^*$  at the gain-maximizing Type 1 criterion 540  
 $k_p$ , with an adjustment for conservatism ( $\Omega_{2,gain+cons}$ ). By adjusting the confidence criteria 541  
so that the counterfactual Type 1 criterion tracks the actual Type 1 criterion, payoffs are 542  
inappropriately incorporated into confidence judgments. As a consequence, higher relative 543  
reward or cost will make a person more likely to select that alternative and, on average, more 544  
confident about reporting that outcome when they do. In effect, this is a naïve optimism 545  
for selecting the highly rewarded outcome and disproportionate pessimism for selecting the 546  
costly outcome: “this highly rewarding perceptual alternative that I have selected is certainly 547  
the state of the world” or “it is costly to me, so it cannot be true”. This bias for higher 548  
confidence with greater reward value (or smaller loss value) is consistent with what has been 549  
reported previously in the perceptual lottery tasks of Lebreton et al. (2018). Yet, we note 550

that a failure to understand the task instructions could have produced the bias we found. 551  
It is possible that observers did not report the probability they were correct, as per the 552  
experimenter instructions, but instead reported something about the expected gain of the 553  
trial when reflecting on their confidence. 554

An inability to appropriately dissociate Type 1 and Type 2 responses, in both subsets of 555  
observers, is compelling. If this is a true inability for sensory decision-making, then there is 556  
a trade-off between maximizing gains for discrimination and accuracy for confidence. That 557  
is, if observers cannot selectively decouple their  $k_1$  and  $k_1^*$  for asymmetric payoffs, then 558  
perhaps they reach a compromise by sacrificing some gains in the Type 1 task by shifting 559  
 $k_1$  toward the accuracy-maximizing criterion  $k_p$ , thereby shifting  $k_1^*$  in a manner that yields 560  
confidence reports more consistent with the objective probability of being correct. Consider, 561  
for example, judging whether an aircraft is heading for collision with an upcoming mountain 562  
peak. The high cost of collision should bias heading judgments toward predicting a collision, 563  
so corrective actions can be taken. But you wouldn't want to be confident in that judgment 564  
just because it results in high cost, so you reduce the bias and are a bit more confident 565  
you'll pass by unscathed. The ideal trade-off between incorporating the payoff structure 566  
versus accurately and confidently making a decision will of course depend on the decision at 567  
hand. Subsequent laboratory experiments can attempt to shift this trade-off by using more 568  
complex reward structures that incorporate both Type 1 and Type 2 judgments. 569

Finally, we turn to the result that the best-fitting Type 2 model had conservatism ap- 570  
plied to the counterfactual criterion  $k_1^*$ . It is currently a matter of debate whether the same 571  
internal measurement of the sensory event is used by both the perceptual and the metacogni- 572  
tive decision-making systems (e.g., Resulaj et al., 2009; Fleming and Daw, 2017). The SDT 573  
framework used here assumes the same internal measurement is used for both judgments. 574  
The Type 2 judgment is thought to include additional noise (Maniscalco and Lau, 2012; 575  
Fleming and Lau, 2014; Bang et al., 2018), and as such, we incorporated reduced metacogni- 576  
tive sensitivity in our modeling. Additionally, our results suggest several possible scenarios 577  
about how the decision boundaries during Type 1 and Type 2 decisions are related. The 578  
Type 1 and Type 2 processes may be computed jointly using the same information, but there 579  
is considerable evidence that neural processing occurs in distinct regions for perceptual and 580

metacognitive decision-making (Shimamura, 2000; Fleming and Dolan, 2012; Rahnev et al., 2016; Shekhar and Rahnev, 2018). The Type 1 system may convey information to the Type 2 system about its decision boundary, or convey only relative information. Additionally, the processes responsible for conservatism are also applied to the counterfactual criterion in the Type 2 system. Given the complexity of the conservatism we observed, it would appear unlikely for the Type 2 system to recreate the phenomenon of conservatism with information acquired independently from the Type 1 system. Thus, we favor the interpretation that the exact effects of priors and payoffs in perceptual decision-making are also propagated to the metacognitive system. Given that a subset of observers were able to dissociate  $k_1$  and  $k_1^*$  by keeping the latter fixed at the neutral criterion, it is less likely that the Type 2 system receives an internal measurement that is coded relative to the discrimination criterion  $k_1$ . Yet, if this is the case, why weren't observers able to reduce the influence of the payoff structure at the second processing step? Further work is required to understand why optimal metacognitive behavior was not achieved.

### 6.3 Conclusion

By manipulating priors and payoffs in a perceptual task, we found sub-optimal decision making at the Type 1 and Type 2 levels. Discrimination judgments were conservative, with no strong tendency for greater conservatism for payoffs than priors. There was also evidence against additivity of criterion shifts for asymmetric priors and payoffs. Confidence judgments were sub-optimal in one of two ways: 1) observers did not consider the role of priors or 2) they incorporated payoffs. Both of these strategies hinder decision-making. For example, a radiologist who ignores prior probabilities when assigning confidence might hesitate recommending further tests for a patient who is a heavy smoker. Similarly, a radiologist who inappropriately incorporates payoffs may be more confident in a positive diagnosis if he receives kickbacks from the imaging center to encourage future scans. The patterns of behavior found in this task point to explanations of why humans may consider trade-offs between maximizing gain and maximizing accuracy, as well as provide new insights about the role of the decision boundary in Type 1 versus Type 2 computations.

## Author Note

609

This work was supported by NIH grant EY08266, the NIH Training Program in Computational Neuroscience grant T90DA043219, the NSF Collaborative Research in Computational Neuroscience grant 1420262, NSF GRFP DGE1342536, and the French ANR grant ANR-18-CE28-0015-01 "VICONTE". This research was presented at the 2018 Vision Sciences Society meeting at St. Pete Beach, Florida.

## References

615

- Ackermann, J. F. and Landy, M. S. (2015). Suboptimal decision criteria are predicted by subjectively weighted probabilities and rewards. *Attention, Perception, and Psychophysics*, 77(2):638–658.
- Bang, J. W., Shekhar, M., and Rahnev, D. (2018). Sensory noise increases metacognitive efficiency. *Journal of Experimental Psychology: General*.
- Beran, M. J., Brandl, J. L., Perner, J., and Joëlle, P., editors (2012). *Foundations of Metacognition*. Oxford University Press, Oxford, UK.
- Brainard, D. H. (1997). The psychophysics toolbox. *Spatial Vision*, 10:433–436.
- Busemeyer, J. R. and Myung, I. J. (1992). An adaptive approach to human decision making: Learning theory, decision theory, and human performance. *Journal of Experimental Psychology: General*, 121(2):177–194.
- Clarke, F. R., Birdsall, T. G., and Tanner, W. P. (1959). Two types of ROC curves and definitions of parameters. *The Journal of the Acoustical Society of America*, 31(5):629–630.
- Dunning, D., Griffin, D. W., Milojkovic, J. D., and Ross, L. (1990). The overconfidence effect in social prediction. *Journal of Personality and Social Psychology*, 58(4):568–581.
- Fleming, S. M. and Daw, N. D. (2017). Self-evaluation of decision-making: A general Bayesian framework for metacognitive computation. *Psychological Review*, 124(1):91–114.

- Fleming, S. M. and Dolan, R. J. (2012). Neural basis of metacognition. *Philosophical Transactions of the Royal Society B Biological Sciences*, 367(1594):1338–1349.
- Fleming, S. M. and Lau, H. C. (2014). How to measure metacognition. *Frontiers in Human Neuroscience*, 8:1–9.
- Green, D. M. and Swets, J. A. (1966). *Signal Detection Theory and Psychophysics*. Wiley, New York.
- Healy, A. F. and Kubovy, M. (1978). The effects of payoffs and prior probabilities on indices of performance and cutoff location in recognition memory. *Memory & Cognition*, 6(5):544–553.
- Healy, A. F. and Kubovy, M. (1981). Probability matching and the formation of conservative decision rules in a numerical analog of signal detection. *Journal of Experimental Psychology: Human Learning and Memory*, 7(5):344–354.
- Horowitz, T. S. (2017). Prevalence in visual search: From the clinic to the lab and back again. *Japanese Psychological Research*, 59(2):65–108.
- Kiani, R. and Shadlen, M. N. (2009). Representation of confidence associated with a decision by neurons in the parietal cortex. *Science*, 324(5928):759–764.
- Kleiner, M., Brainard, D., Pelli, D., Ingling, A., Murray, R., and Broussard, C. (2007). What’s new in psychtoolbox-3? *Perception*, 36(14):1–16.
- Kubovy, M. (1977). A possible basis for conservatism in signal detection and probabilistic categorization tasks. *Perception & Psychophysics*, 22(3):277–281.
- Lebreton, M., Langdon, S., Slieker, M. J., Nooitgedacht, J. S., Goudriaan, A. E., Denys, D., van Holst, R. J., and Luigjes, J. (2018). Two sides of the same coin: Monetary incentives concurrently improve and bias confidence judgments. *Science Advances*, 4(5):eaq0668.
- Lee, W. and Zentall, T. R. (1966). Factorial effects in the categorization of externally distributed stimulus samples. *Perception & Psychophysics*, 1(2):120–124.

- Macmillan, N. A. and Creelman, C. D. (2005). *Detection Theory: A User's Guide*. Lawrence Erlbaum Associates, Inc., Mahwah, NJ, 2nd edition.
- Maddox, W. T. (2002). Toward a unified theory of decision criterion learning in perceptual categorization. *Journal of the Experimental Analysis of Behavior*, 78(3):567–595.
- Maddox, W. T. and Bohil, C. J. (1998). Base-rate and payoff effects in multidimensional perceptual categorization. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 24(6):1459–1482.
- Maddox, W. T. and Bohil, C. J. (2000). Costs and benefits in perceptual categorization. *Memory & Cognition*, 28(4):597–615.
- Maddox, W. T. and Bohil, C. J. (2003). A theoretical framework for understanding the effects of simultaneous base-rate and payoff manipulations on decision criterion learning in perceptual categorization. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 29(2):307–320.
- Maddox, W. T. and Dodd, J. L. (2001). On the relation between base-rate and cost-benefit learning in simulated medical diagnosis. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 27(6):1367–1384.
- Mamassian, P. (2008). Overconfidence in an objective anticipatory motor task. *Psychological Science*, 19(6):601–606.
- Mamassian, P. (2016). Visual confidence. *Annual Review of Vision Science*, 2:459–481.
- Manis, M., Dovalina, I., Avis, N. E., and Cardoze, S. (1980). Base rates can affect individual predictions. *Journal of Personality and Social Psychology*, 38(2):231–248.
- Maniscalco, B. and Lau, H. C. (2012). A signal detection theoretic approach for estimating metacognitive sensitivity from confidence ratings. *Consciousness and Cognition*, 21:422–430.
- Maniscalco, B. and Lau, H. C. (2014). Signal detection theory analysis of type 1 and type 2 data: meta- $d'$ , response-specific meta- $d'$ , and the unequal variance SDT model. In

- Fleming, S. M. and Frith, C. D., editors, *The Cognitive Neuroscience of Metacognition*, 684  
pages 25–66. Springer. 685
- Metcalfe, J. and Shimamura, A. P., editors (1996). *Metacognition: Knowing about Knowing*. 686  
MIT Press, Cambridge, MA. 687
- Pelli, D. G. (1997). The VideoToolbox software for visual psychophysics: Transforming 688  
numbers into movies. *Spatial Vision*, 10(4):437–442. 689
- Pouget, A., Drugowitsch, J., and Kepecs, A. (2016). Confidence and certainty: distinct 690  
probabilistic quantities for different goals. *Nature Neuroscience*, 19(3):366–374. 691
- Rahnev, D., Nee, D. E., Riddle, J., Larson, A. S., and D’Esposito, M. (2016). Causal 692  
evidence for frontal cortex organization for perceptual decision making. *Proceedings of the* 693  
*National Academy of Sciences of the United States of America*, 113(21):6059–6064. 694
- Resulaj, A., Kiani, R., Wolpert, D. M., and Shadlen, M. N. (2009). Changes of mind in 695  
decision-making. *Nature*, 461:263–266. 696
- Rigoux, L., Stephan, K. E., Friston, K. J., and Daunizeau, J. (2014). Bayesian model 697  
selection for group studies—Revisited. *NeuroImage*, 84:971–985. 698
- Shekhar, M. and Rahnev, D. (2018). Distinguishing the Roles of Dorsolateral and Anterior 699  
PFC in Visual Metacognition. *The Journal of Neuroscience*, 38(22):5078–5087. 700
- Sherman, M. T., Seth, A. K., Barrett, A. B., and Kanai, R. (2015). Prior expectations 701  
facilitate metacognition for perceptual decision. *Consciousness and Cognition*, 35:53–65. 702
- Shimamura, A. P. (2000). Toward a cognitive neuroscience of metacognition. *Consciousness* 703  
*and Cognition*, 9:313–323. 704
- Smith, J. D., Shields, W. E., and Washburn, D. A. (2003). The comparative psychology of 705  
uncertainty monitoring and metacognition. *Behavioral and Brain Sciences*, 26(3):317–339. 706
- Stephan, K. E., Penny, W. D., Daunizeau, J., Moran, R. J., and Friston, K. J. (2009). 707  
Bayesian model selection for group studies. *NeuroImage*, 46(4):1004–1017. 708



- Stevenson, M. K., Busemeyer, J. R., and Naylor, J. C. (1990). Judgement and decision- 709  
making theory. In Dunnette, M. D. and Hough, L. M., editors, *Handbook of Industrial* 710  
*and Organizational Psychology*, pages 283–374. Consulting Psychologists Press, Palo Alto, 711  
CA. 712
- Ulehla, Z. J. (1966). Optimality of perceptual decision criteria. *Journal of Experimental* 713  
*Psychology*, 71(4):564–569. 714
- Wolfe, J. M., Horowitz, T. S., and Kenner, N. M. (2005). Rare items often missed in visual 715  
searches. *Nature*, 435:439–440. 716
- Zylberberg, A., Wolpert, D. M., and Shadlen, M. N. (2018). Counterfactual reasoning 717  
underlies the learning of priors in decision making. *Neuron*, 99(5):1083–1097. 718