

A confirmation bias in perceptual decision-making due to hierarchical approximate inference

Richard D. Lange^{1,2,*}, Ankani Chattoraj¹,
Jeffrey M. Beck³, Jacob L. Yates¹, Ralf M. Haefner^{1,*}

¹Brain and Cognitive Sciences, University of Rochester, Rochester, NY 14627, USA.

²Computer Science, University of Rochester, Rochester, NY 14627, USA.

³Department of Neurobiology, Duke University, Durham, NC 27708, USA.

*Corresponding authors: rlange@ur.rochester.edu, rhaefne2@ur.rochester.edu.

January 20, 2020

1 Summary

2 Human decisions are known to be systematically biased. A prominent example of such a bias
3 occurs when integrating a sequence of sensory evidence over time. Previous empirical studies differ
4 in the nature of the bias they observe, ranging from favoring early evidence (primacy), to favoring
5 late evidence (recency). Here, we present a unifying framework that explains these biases and
6 makes novel psychophysical and neurophysiological predictions. By explicitly modeling both the
7 approximate and the hierarchical nature of inference in the brain, we show that temporal biases
8 depend on the balance between “sensory information” and “category information” in the stimulus.
9 Finally, we present new data from a human psychophysics task that confirm that temporal biases
10 can be robustly changed within subjects as predicted by our models.

11 Introduction

12 Imagine a doctor trying to infer the cause of a patient’s symptoms from an x-ray image. Unsure
13 about the evidence in the image, she asks a radiologist for a second opinion. If she tells the
14 radiologist her suspicion, she may bias his report. If she does not, he may not detect a faint
15 diagnostic pattern. As a result, if the evidence in the image is hard to detect or ambiguous,
16 the radiologist’s second opinion, and hence the final diagnosis, may be swayed by the doctor’s
17 initial hypothesis. The problem faced by these doctors exemplifies the difficulty of *hierarchical*
18 *inference*: each doctor’s suspicion both informs and is informed by their collective diagnosis. If
19 they are not careful, their diagnosis may fall prey to circular reasoning. The brain faces a similar
20 problem during perceptual decision-making: any decision-making area combines sequential signals
21 from sensory brain areas, not directly from sensory input, just as the doctors’ consensus is based
22 on their individual diagnoses rather than on the evidence *per se*. If sensory signals in the brain
23 themselves reflect inferences that combine both prior expectations and sensory evidence, we suggest
24 that this can then lead to an observable *perceptual* confirmation bias (Nickerson, 1998).

25 We formalize this idea in the context of approximate Bayesian inference and classic evidence-
26 integration tasks in which a range of biases has been observed and for which a unifying explanation

27 is currently lacking. Evidence-integration tasks require subjects to categorize a sequence of inde-
28 pendent and identically distributed (iid) draws of stimuli (Gold and Shadlen, 2007; Bogacz et al.,
29 2006). Previous normative models of evidence integration hinge on two quantities: the amount of
30 information available on a single stimulus draw and the total number of draws. One might expect,
31 then, that temporal biases should have some canonical form in tasks where these quantities are
32 matched. However, existing studies are heterogeneous, reporting one of three distinct motifs: some
33 find that early evidence is weighted more strongly (a primacy effect) (Kiani et al., 2008; Nienborg
34 and Cumming, 2009) some that information is weighted equally over time (as would be optimal)
35 (Wyart et al., 2012; Brunton et al., 2013; Raposo et al., 2014), and some find late evidence being
36 weighted most heavily (a recency effect) (Drugowitsch et al., 2016) (Figure 1a,c). While there
37 are myriad differences between these studies such as subject species, sensory modality, stimulus
38 parameters, and computational frameworks (Kiani et al., 2008; Brunton et al., 2013; Glaze et al.,
39 2015; Drugowitsch et al., 2016), none of these aspects alone can explain their different findings.

40 We extend classic evidence-integration models to the *hierarchical* case by including an explicit
41 intermediate sensory representation, analogous to modeling each doctor’s individual diagnosis in
42 addition to their consensus in the example above (Figure 1b). Taking this intermediate inference
43 stage into account makes explicit that task difficulty is modulated by two distinct types of informa-
44 tion exposing systematic differences between existing tasks: the information between the stimulus
45 and sensory representation (“sensory information”), and the information between sensory represen-
46 tation and category (“category information”) (Figure 1b). These differences alone do not entail any
47 bias as long as inference is exact. However, inference in the brain is necessarily *approximate* and
48 this approximation can interfere with its ability to account for its own biases. Implementing two
49 approximate hierarchical inference algorithms, we find that they both result in biases in agreement
50 with our data, and can indeed explain the puzzling discrepancies in the literature.

51 Results

52 “Sensory Information” vs “Category Information”

Normative models of decision-making in the brain are typically based on the idea of an *ideal observer*, who uses Bayes’ rule to infer the most likely category on each trial given the stimulus. On each trial in a typical task, the stimulus consists of multiple “frames” presented in rapid succession. (By “frames” we refer to discrete independent draws of stimulus values that are not necessarily visual). If the evidence in each frame, e_f , is independent then evidence can be combined by simply multiplying the associated likelihoods. And if the categorical identity of the stimulus is a binary variable $C \in \{-1, +1\}$, then this process corresponds to the famous sequential ratio test summing the log odds implied by each piece of evidence (Wald and Wolfowitz, 1948; Bogacz et al., 2006):

$$\begin{aligned} p(C = +1|e_1, \dots, e_F) &\propto p(C = +1) \prod_{f=1}^F p(e_f|C = +1) \\ \log \frac{p(C = +1|e_1, \dots, e_F)}{p(C = -1|e_1, \dots, e_F)} &= \log \frac{p(C = +1)}{p(C = -1)} + \sum_{f=1}^F \log \frac{p(e_f|C = +1)}{p(e_f|C = -1)}. \end{aligned}$$

53 As a result, the ideal observer’s performance is determined by (i) the information about C available
54 on each frame, $p(e_f|C)$, and (ii) the number of frames per trial.

55 However, in the brain, any decision-making area does not base its decision on the externally
56 presented stimulus directly, but rather on an intermediate sensory representation of the stimulus.

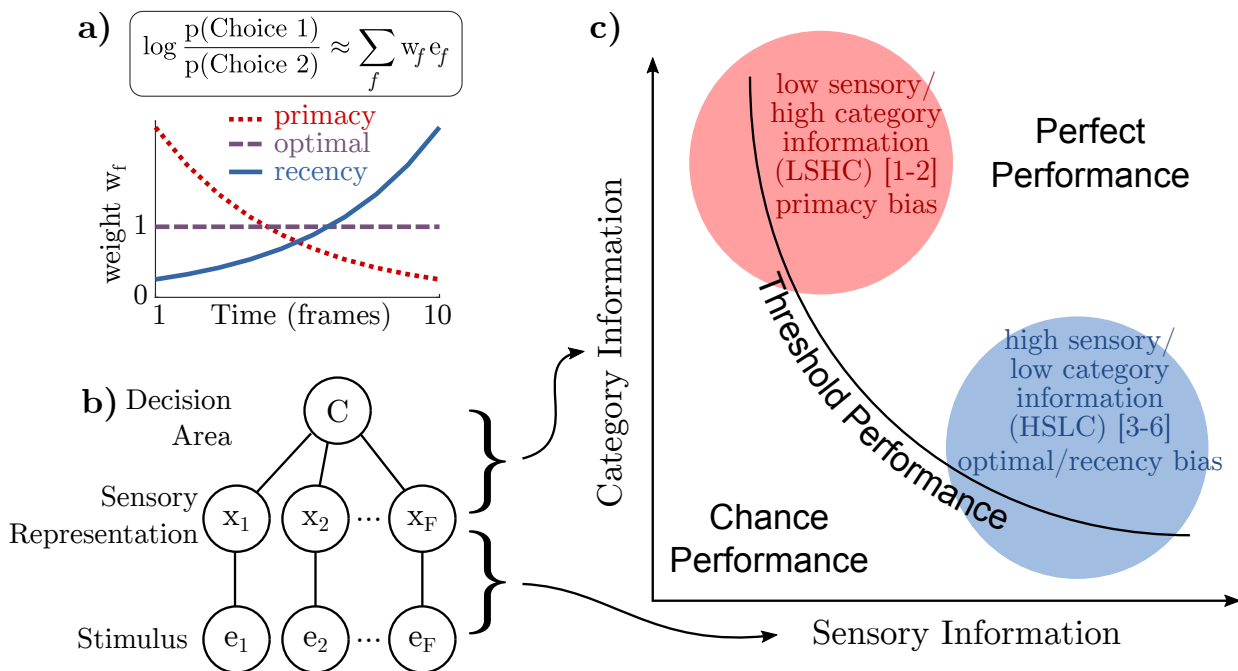


Figure 1: **a)** A subject’s “temporal weighting strategy” is an estimate of how their choice is based on a weighted sum of each frame of evidence e_f . Three commonly observed motifs are decreasing weights (primacy), constant weights (optimal), or increasing weights (recency). **b)** Information in the stimulus about the category may be decomposed into information in each frame about a sensory variable (“sensory information”) and information about the category given the sensory variable (“category information”). **c)** Category information and sensory information may be manipulated independently, creating a two-dimensional space of possible tasks. Any level of task performance can be the result of different combinations of sensory and category information. A qualitative placement of previous work into this space separates those that find primacy effects in the upper-left from those that find recency effects or optimal weights in the lower right (see Supplemental Text for detailed justification). Numbered references are: [1] Kiani et al., [2] Nienborg and Cumming, [3] Brunton et al., [4] Wyart et al., [5] Raposo et al., [6] Drugowitsch et al.

57 This intermediate representation is itself often assumed to be the result of an inference process in which
58 sensory neurons compute the posterior distribution $p(\mathbf{x}|e)$ (Fiser et al., 2010; Pouget et al., 2013;
59 Gershman and Beck, 2016; Lange and Haefner, 2020) over some latent variable \mathbf{x} given the external
60 evidence e in an internal model of the world (Mumford, 1992; Lee and Mumford, 2003; Yuille
61 and Kersten, 2006). This process is naturally formalized as hierarchical inference (Figure 1b).
62 This implies that the information between the stimulus and category can be partitioned into the
63 information between the stimulus and the sensory representation (e to \mathbf{x}), and the information
64 between sensory representation and category (\mathbf{x} to C). We call these “sensory information” and
65 “category information,” respectively (Figure 1b). These two kinds of information define a two-
66 dimensional space in which a given task is located as a single point (Figure 1c). For example, in
67 a visual task each e_f would be the image on the screen while x_f might be image patches that are
68 assumed to be sparsely combined to form the image (Olshausen and Field, 1997). The posterior
69 over the latent features x_f would be represented by the activity of relevant neurons in visual cortex.

70 An evidence integration task may be challenging either because each frame is perceptually
71 unclear (low “sensory information”), or because each frame alone is insufficient to determine the
72 category for the whole trial (low “category information”). Consider the classic dot motion task
73 (Newsome and Pare, 1988) and the Poisson clicks task (Brunton et al., 2013), which occupy opposite
74 locations in the space. In the classic low-coherence dot motion task, subjects view a cloud of moving
75 dots, a small percentage of which move “coherently” in one direction. Here, sensory information
76 is low since the percept of net motion is weak on each frame. Category information, on the other
77 hand, is high, since knowing the true net motion on a single frame would be highly predictive of
78 the correct choice (and of motion on subsequent frames). In the Poisson clicks task on the other
79 hand, subjects hear a random sequence of clicks in each ear and must report the side with the
80 higher rate. Here, sensory information is high since each click is well above sensory thresholds.
81 Category information, however, is low, since knowing the side on which a single click was presented
82 provides only little information about the correct choice for the trial as a whole (and the side of the
83 other clicks). Another way to think about category information is as “temporal coherence” of the
84 stimulus: the more each frame of evidence is predictive of the correct choice, the more the frames
85 must be predictive of each other, whether a frame consists of visual dots or of auditory clicks. Note
86 that our distinction between sensory and category information is different from the well-studied
87 distinction between internal and external noise; in general, both internal and external noise will
88 reduce the amount of sensory and category information.

89 Optimal inference requires accounting for all possible sources of information. Ideally, then,
90 sensory areas would not only represent the current evidence, $p(x_f|e_f)$, but should incorporate prior
91 information based on previous frames to compute $p(x_f|e_1, \dots, e_f)$. While the sensory area no longer
92 has direct access to the earlier frames, this is mathematically equivalent to using the current belief
93 in the category C as a prior:

$$p(x_f|e_1, \dots, e_f) \propto p(e_f|x_f) \sum_c \underbrace{p(C=c|e_1, \dots, e_{f-1})}_{p_{f-1}(C)} p(x_f|C=c). \quad (1)$$

94 Mechanistically, this suggests that the brain’s running estimate of the category, $p_{f-1}(C)$, should
95 be continuously fed back to sensory areas, acting as a prior that biases the representation to agree
96 with the current belief about the category (Lee and Mumford, 2003; Haefner et al., 2016; Tajima
97 et al., 2016; Lange and Haefner, 2020). Importantly, such a bias is optimal in the sense that it
98 makes instantaneous sensory estimates more accurate. Despite this instantaneous sensory bias,
99 *exact* inference in this model does *not* induce any bias in the posterior over the category C . That
100 is, although the ideal observer’s inference about x_f is biased by e_1, \dots, e_{f-1} , this bias is removed

101 by precisely accounting for it in the update to $p_f(C)$ (Zylberberg et al., 2018).

102 Unlike the ideal observer, inference in the brain is necessarily *approximate* ((Fiser et al.,
103 2010; Pouget et al., 2013) and the implications of this fact on evidence integration has so far been
104 unknown. Below, we consider two models, each implementing approximate hierarchical inference in
105 one of the two major classes of approximate inference schemes known from statistics and machine
106 learning: sampling-based and variational inference (Bishop, 2006; Murphy, 2012), both of which
107 have been previously proposed models for neural inference (Fiser et al., 2010; Pouget et al., 2013).
108 In both models, a confirmation bias arises as a direct consequence of the approximate nature
109 of inference over the intermediate sensory variables in the brain. The strength of the predicted
110 confirmation bias depends directly on the amount of category information in the stimulus, since that
111 governs how strongly past frames inform inferences about the present frame. Our models predict an
112 overweighting of early evidence when sensory information is low and category information is high,
113 but not when sensory information is high and category information is low, even when performance
114 is matched in both conditions (Fig. 1c, for model details see “Approximate inference models”
115 section below).

116 Qualitatively placing prior studies in the space spanned by these two kinds of information results
117 in two clusters: the studies that report primacy effects are located in the upper left quadrant (low-
118 sensory/high-category or LSHC) and studies with flat weighting or recency effects are in the lower
119 right quadrant (high-sensory/low-category or HSLC) (Figure 1c). This initially suggests that the
120 trade-off between sensory information and category information may indeed underlie differences in
121 temporal weighting seen in previous studies. Further, this framework allows us to make new and
122 easily testable predictions for how simple changes in stimulus statistics of previous studies should
123 change the temporal weighting they find (Supplemental Table S1). We next describe a novel set
124 of visual discrimination tasks designed to directly probe this trade-off between sensory information
125 and category information to test these predictions within individual subjects.

126 Visual Discrimination Task

127 We designed a visual discrimination task with two stimulus conditions that correspond to the two
128 opposite sides of this task space, while keeping all other aspects of the design the same (Figure 2a).
129 If our theory is correct, then we should be able to change individual subjects’ temporal weighting
130 strategy simply by changing the sensory-category information trade-off.

131 The stimulus in our task consisted of a sequence of ten visual frames (83ms each). Each frame
132 consisted of band-pass-filtered white noise with excess orientation power either in the -45° or the
133 $+45^\circ$ orientation (Beaudot and Mullen, 2006) (Figure 2b,d). On each trial, there was a single true
134 orientation category, but individual frames might differ in their orientation. At the end of each
135 trial, subjects reported whether the stimulus was oriented predominantly in the -45° or the $+45^\circ$
136 orientation. The stimulus was presented as an annulus around the fixation marker in order to
137 minimize the effect of small fixational eye movements (Methods).

138 If the brain’s intermediate sensory representation reflects the orientation in each frame, then
139 sensory information in our task is determined by how well each frame determines the orientation
140 of that frame (i.e. the amount of “noise” in each frame), and category information is determined
141 by the probability that any given frame’s orientation matches the trial’s category. We chose to
142 quantify both sensory information and category information, using signal detection theory, as the
143 area under the receiver-operating-characteristic curve for e_f and x_f (sensory information), or for x_f
144 and C (category information). Hence for a ratio of 5 : 5, a frame’s orientation does not predict the
145 correct choice and category information is 0.5. For a ratio of 10 : 0, knowledge of the orientation of

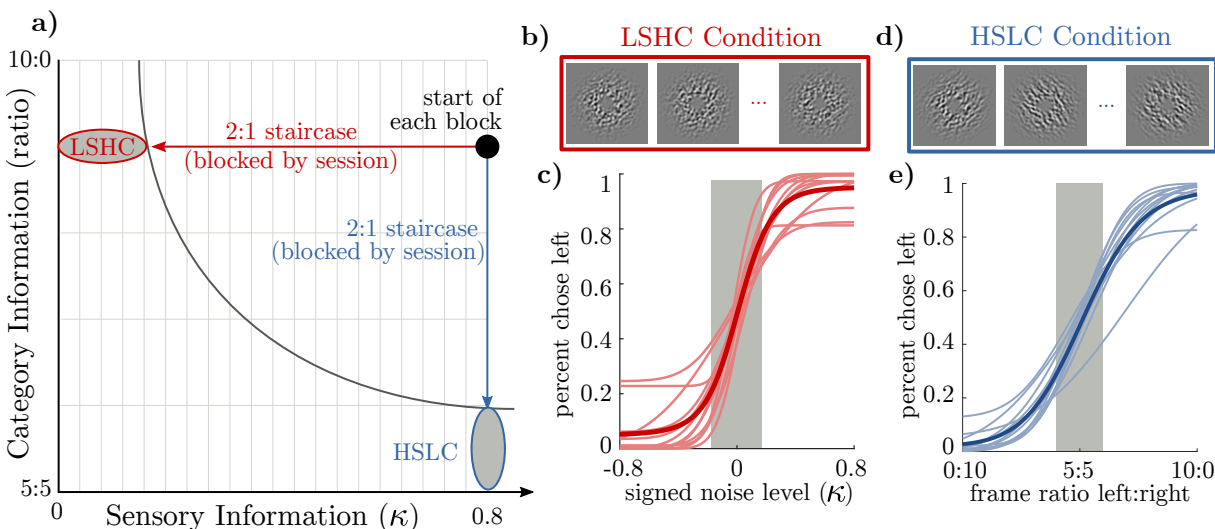


Figure 2: Summary of experiment design. **a)** Category information is determined by the expected ratio of frames in which the orientation matches the correct category, and sensory information is determined by a parameter κ determining the degree of spatial orientation coherence (Methods). At the start of each block, we reset the staircase to the same point, with category information at 9 : 1 and κ at 0.8. We then ran a 2-to-1 staircase either on κ or on category information. The LSHC and HSLC ovals indicate sub-threshold trials; only these trials were used in the regression to infer subjects' temporal weights. **b)** Visualization of a noisy stimulus in the LSHC condition. All frames are oriented to the right. **c)** Psychometric curves for all subjects (thin lines) and averaged (thick line) over the κ staircase. Shaded gray area indicates the median threshold level across all subjects. **d)** Example frames in the HSLC condition. The orientation of each frame is clear, but orientations change from frame to frame. **e)** Psychometric curves over frame ratios, plotted as in (c).

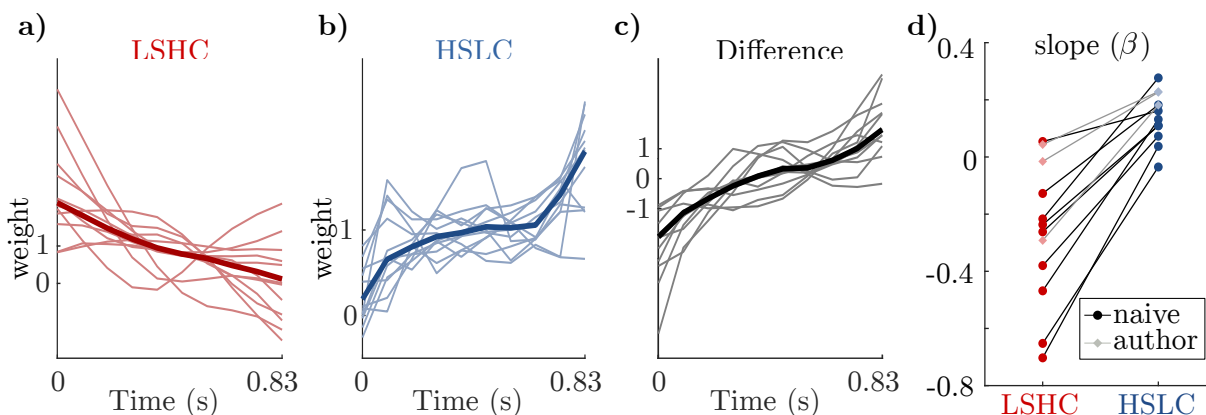


Figure 3: Subjects' temporal weights. **a-b)** Temporal weights for individual subjects (thin lines) and the mean across all subjects (thick lines). Weights are normalized to have a mean of 1 to emphasize shape rather than magnitude. Individual subjects' curves were fit using a cross-validated smoothness term for visualization purposes only (Methods). **c)** Difference of normalized weights (HSLC–LSHC). Despite variability across subjects in (a-b), each subject reliably changes in the direction of a recency effect. **d)** *Change* in slope between the two task contexts for each subject is consistently positive. Points are median slope values after bootstrap-resampling of the data. We summarize subjects' temporal weighting strategy with an exponential fit; the slope parameter $\beta > 0$ corresponds to recency and $\beta < 0$ to primacy (similar results for linear fits, see SI).

146 a single frame is sufficient to determine the correct choice and category information is 1. Exactly
147 quantifying sensory information depends on individual subjects, but likewise ranges from 0.5 to 1.
148 For a more detailed discussion, see Supplementary Text.

149 Using this stimulus, we tested 12 human subjects (9 naive and 3 authors) comparing two
150 conditions intended to probe the difference between the LSHC and HSLC regimes. Starting with
151 both high sensory and high category information, we either ran a 2:1 staircase lowering the sensory
152 information while keeping category information high, or we ran a 2:1 staircase lowering category
153 information while keeping sensory information high (Figure 2a). These are the LSHC and HSLC
154 conditions, respectively (Figure 2b,d). For each condition, we used logistic regression to infer,
155 for each subject, the influence of each frame onto their choice. Subjects' overall performance was
156 matched in the two conditions by setting a performance threshold below which trials were included
157 in the analysis (Methods).

158 In agreement with our hypothesis, we find predominantly flat or decreasing temporal weights
159 in the LSHC condition (Figure 3a). However, when the information is partitioned differently –
160 in the HSLC condition – we find flat or increasing weights (Figure 3b). Importantly, despite
161 variability between subjects in each condition, a within-subject comparison revealed that the change
162 in slope between the two conditions was as predicted for all subjects (Figure 2c,d) ($p < 0.05$ for
163 10 of 12 subjects, bootstrap). This demonstrates that the trade-off between sensory and category
164 information in a task robustly changes subjects' temporal weighting strategy as we predicted, and
165 further suggests that the sensory-category information trade-off may resolve the discrepant results
166 in the literature.

167 Approximate inference models

168 We will now show that these significant changes in evidence weighting for different stimulus statis-
 169 tics arise naturally in common models of how the brain might implement approximate inference.
 170 In particular, we show that both a neural sampling-based approximation (Hoyer and Hyvärinen,
 171 2003; Fiser et al., 2010; Haefner et al., 2016; Orbán et al., 2016) and a parametric (mean-field)
 172 approximation (Beck et al., 2013; Raju and Pitkow, 2016) can explain the observed pattern of
 173 changing temporal weights as a function of stimulus statistics.

174 Optimal inference in our task, as in other evidence integration tasks, requires computing the
 175 posterior over C conditioned on the evidence e_1, \dots, e_f , which can be expressed as the Log Posterior
 176 Odds (LPO),

$$\underbrace{\log \frac{p(C = +1|e_1, \dots, e_f)}{p(C = -1|e_1, \dots, e_f)}}_{\text{LPO}_f} = \log \frac{p(C = +1)}{p(C = -1)} + \sum_{i=1}^f \underbrace{\log \frac{p(e_i|C = +1)}{p(e_i|C = -1)}}_{\text{LLO}_i}, \quad (2)$$

177 where LLO_f is the log likelihood odds for frame f (Gold and Shadlen, 2007; Bogacz et al., 2006).
 178 To reflect the fact that the brain has access to only one frame of evidence at a time, this can
 179 be rewritten this as an *online* update rule, summing the previous frame’s log posterior with new
 180 evidence gleaned on the current frame:

$$\text{LPO}_f = \text{LPO}_{f-1} + \text{LLO}_f. \quad (3)$$

181 This expression is derived from the ideal observer and is still exact. Since the ideal observer weights
 182 all frames equally, the *online* nature of inference in the brain cannot by itself explain temporal
 183 biases. Furthermore, because performance is matched in the two conditions of our experiment,
 184 their differences cannot be explained by the total amount of information, governed by the likelihood
 185 $p(e_f|C)$.

186 To understand how biases arise, we must examine the log likelihood odds term, LLO, in detail.
 187 In a hierarchical model, computing $p(e_f|C)$ for each C requires marginalizing over the intervening
 188 x_f as follows:

$$\begin{aligned} p(e_f|C) &\propto \int p(e_f|x_f)p(x_f|C)dx_f \\ &\propto \mathbb{E}_{p(x_f|e_f)} \left[\frac{p(x_f|C)}{p(x_f)} \right], \end{aligned} \quad (4)$$

189 This suggests that evidence about the current frame is formed in a two step process: first, x_f is
 190 inferred given e_f , and second an expectation is taken with respect to $p(x_f|e_f)$, where the operand
 191 of the expectation depends only on the relation between x_f and C . No sub-optimality nor biases
 192 have been introduced yet.

193 A key assumption in our models that gives rise to temporal biases is that sensory areas represent
 194 the approximate *posterior* belief over x_f given all available information, i.e. including the earlier
 195 frames in the trial (equation (1)). This assumption differs from some models of inference in the
 196 brain that assume populations of sensory neurons strictly encode the *likelihood* of the stimulus (or
 197 instantaneous posterior) (Ma et al., 2006; Beck et al., 2008), but is consistent with other models
 198 from both sampling and parametric families (Berkes et al., 2011; Haefner et al., 2016; Raju and
 199 Pitkow, 2016; Tajima et al., 2016).

200 As introduced in equation (1), representing the full posterior over x_f implies taking into account
201 all previous frames. In other words, the brain’s belief about x_f depends both on the external
202 evidence, e_f , via the likelihood, but also on the brain’s current belief about C , via the prior. If
203 this were not the case – if sensory areas represented only the instantaneous evidence $p(e_f|x_f)$ –
204 then integrating evidence in an unbiased way would simply be a matter of applying equation (4).
205 However, such an inference scheme comes at the expense of a worse instantaneous representation
206 (Zylberberg et al., 2018).

207 There is thus tension between inferences at two timescales. Instantaneously, it seems advanta-
208 geous to represent $p(x_f|e_1, \dots, e_f)$, while integrating evidence online requires an expectation taken
209 with respect to $p(x_f|e_f)$ (equation (4)). Assuming that the former is represented by sensory areas,
210 the decision area of an approximate ideal observer now needs to correct for, or “subtract out”
211 its influence on those sensory responses. Approximations either to the posterior $p(x_f|e_1, \dots, e_f)$
212 itself or to the bias-correction may underlie the observed behavioral biases. To test this, we imple-
213 mented approximate hierarchical online inference (where “online” means observing a single frame
214 at a time) for a discrimination task using two previously proposed frameworks for how inference
215 might be implemented in neural circuits: neural sampling (Hoyer and Hyvärinen, 2003; Fiser et al.,
216 2010; Haefner et al., 2016; Orbán et al., 2016) and mean field variational inference (Beck et al.,
217 2013; Raju and Pitkow, 2016) (Figure 4).

218 Sampling model

219 The neural sampling hypothesis states that variable neural activity over brief time periods can be
220 interpreted as a sequence of samples from the brain’s posterior over latent variables in its internal
221 model. In our model, samples of x_f are drawn from the full posterior having incorporated the
222 running estimate of $p_f(C)$ (equation (1), Methods), but from equation (4) we would like to use
223 these samples to compute an expectation with respect to only the instantaneous evidence, $p(x_f|e_f)$.
224 The canonical way to compute an expectation with respect to one distribution using samples from
225 another is “importance sampling,” which weights each sample so as to adjust for the difference
226 between the two distributions (Shi and Griffiths, 2009; Murphy, 2012, Chapter 23). In the most
227 extreme case of continual online updates, one could imagine that the brain computes each update
228 to $p_f(C)$ after observing a single sample of x_f . In this case, no correction would be possible; a
229 downstream area would be unable to recover the instantaneous distribution $p(x_f|e_f)$ from a sample
230 sample from the full posterior $p(x_f|e_1, \dots, e_f)$. If the brain is able to base each update on multiple
231 samples, then the *importance weights* of each sample in the update account for the discrepancy
232 between $p(x_f|e_f)$ and $p(x_f|e_1, \dots, e_f)$ (Methods). While this approach is unbiased in the limit of
233 infinitely many samples, it incurs a bias for a finite number – the relevant regime for the brain
234 (Owen, 2013). The bias is *as if* the expectation in (4) is taken with respect to an intermediate
235 distribution that lies between the fully biased one ($p(x_f|e_1, \dots, e_f)$) and the unbiased one ($p(x_f|e_f)$)
236 (Cremer et al., 2017).

237 Under-correcting for the prior that was fed back results in a positive feedback loop between
238 decision-making and sensory areas which we call a “perceptual confirmation bias.” Importantly,
239 this feedback loop is strongest when category information is high, corresponding to stronger feed-
240 back, and sensory information is low, since that makes x_f less dependent on e_f . Figure 4b and
241 Supplemental Figure S5a-c show performance for the ideal observer and for the resulting sampling-
242 based model, respectively, across all combinations of sensory and category information. White lines
243 show threshold performance (70% correct) as in Figure 1c.

244 This model reproduces the primacy effect, and how the temporal weighting changes as the
245 stimulus information changes seen in previous studies. Importantly, it predicted the same within-

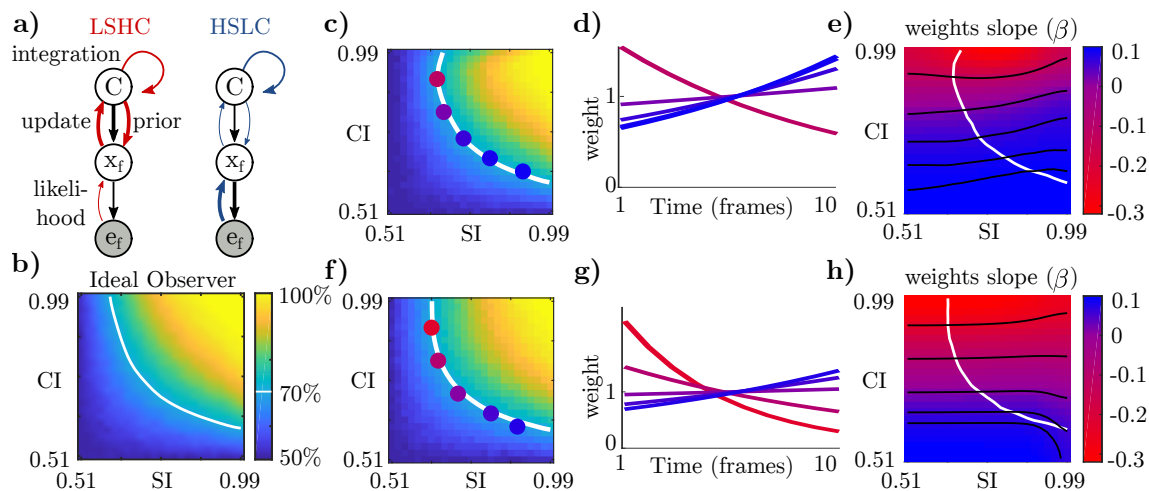


Figure 4: Approximate inference models explain results. **a)** The difference in stimulus statistics between HSLC and LSHC trade-offs implies that the relevant sensory representation is differentially influenced by the stimulus or by beliefs about the category C . A “confirmation bias” or feedback loop between x and C emerges in the LSHC condition but is mitigated in the HSLC condition. Black lines indicate the underlying generative model, and red/blue lines indicate information flow during inference. Arrow width represents coupling strength. **b)** Performance of an ideal observer reporting C given ten frames of evidence. White line shows threshold performance, defined as 70% correct. **c)** Performance of the sampling model with $\gamma = 0.1$. Colored dots correspond to lines in the next panel. **d)** Temporal weights in the model transition from recency to a strong primacy effect, all at threshold performance, as the stimulus transitions from the high-sensory/low-category to the low-sensory/high-category conditions. **e)** Using the same exponential fit as used with human subjects, visualizing how temporal biases change across the entire task space. Red corresponds to primacy, and blue to recency. White contour as in (c). Black lines are iso-contours for slopes corresponding to highlighted points in (c). **f-h)** Same as **c-d** but for the variational model with $\gamma = 0.1$.

246 subject change seen in our data (Haefner et al., 2016). However, double-counting the prior alone
247 cannot explain recency effects (Supplemental Figure S5a-c,j-l).

248 There are two simple and biologically-plausible explanations for the observed recency effect
249 which turn out to be mathematically equivalent. First, the brain may try to actively compensate
250 for the prior influence on the sensory representation by subtracting out an estimate of that influence.
251 That is, the brain could do approximate bias correction to mitigate the effect of the confirmation
252 bias. We modeled linear bias correction by explicitly subtracting out a fraction of the running
253 posterior odds at each step:

$$\text{LPO}_f \leftarrow \text{LPO}_{f-1}(1 - \gamma) + \hat{\text{LLO}}_f \quad (5)$$

254 where $0 \leq \gamma \leq 1$ and $\hat{\text{LLO}}_f$ is the model's (biased) estimate of the log likelihood odds. Second, the
255 brain may assume a non-stationary environment, i.e. C is not constant over a trial. Interestingly,
256 Glaze et al. (2015) showed that optimal inference in this case implies equation (5), which can
257 be interpreted as a noiseless, discrete time version of the classic drift-diffusion model (Gold and
258 Shadlen, 2007) with γ as a leak parameter.

259 Incorporating equation (5) into our model reduces the primacy effect in the upper left of the
260 task space and leads to a recency effect in the lower right (Figure 4c-e, Supplemental Figure S5),
261 as seen in the data.

262 Variational model

263 The second major class of models for how probabilistic inference may be implemented in the brain
264 – based on mean-field parametric representations (Ma et al., 2006; Beck et al., 2013) – behaves
265 similarly. These models commonly assume that distributions are encoded *parametrically* in the
266 brain, but that the brain explicitly accounts for dependencies only between subsets of variables, e.g.
267 within the same cortical area. (Raju and Pitkow, 2016). We therefore make the assumption that
268 the joint posterior $p(x, C|e)$ is approximated in the brain by a product of parametric distributions,
269 $q(x)q(C)$ (Beck et al., 2013; Raju and Pitkow, 2016). Inference proceeds by iteratively minimizing
270 the Kullback-Leibler divergence between $q(x)q(C)$ and $p(x, C|e)$ (Methods). As in the sampling
271 model, the current belief about the category C acts as a prior over x . Because this model is unable
272 to explicitly represent posterior dependencies between sensory and decision variables, both x and
273 C being positive and both x and C being negative act as attractors of its temporal dynamics.
274 This yields qualitatively the same behavior as the sampling model: a stronger influence of early
275 evidence and a transition from primacy to flat weights as category information decreases. As in the
276 sampling model, recency effects emerge only when approximate bias correction is added (Figure
277 4f-h, Supplemental Figure S5j-r). Whereas the limited number of samples was the key deviation
278 from optimality in the sampling model, here it is the assumption that the brain represents its beliefs
279 separately about x and C in a factorized form (Methods).

280 Optimal bias correction

281 A leak term implements optimal inference in a changing environment (Glaze et al., 2015), but each
282 trial of our task is stationary. One might therefore expect that a leak term, or $\gamma > 0$, would impair
283 the model's performance in our task. On the other hand, we motivated the leak term by suggesting
284 that it could approximately correct for the confirmation bias. Under this second interpretation,
285 one might instead expect performance to *improve* for some $\gamma > 0$, especially for conditions where
286 the confirmation bias was strong.

287 We investigated the relationship between the leak (γ) and model performance. First, we sim-
288 ulated the importance sampling model with $\gamma = 0.1$ and $\gamma = 0.5$ and compared its performance

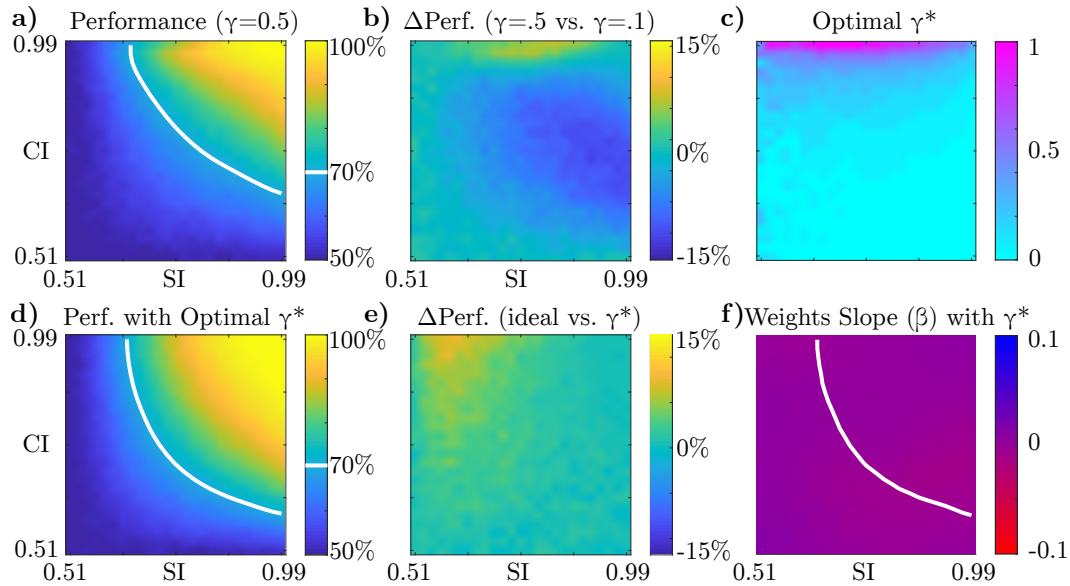


Figure 5: Optimizing performance with respect to γ (see also Supplemental Figure S6). **a)** Model performance across task space with $\gamma = 0.5$ (compare with Figure 4c in which $\gamma = 0.1$). **b)** Difference in performance for $\gamma = 0.5$ versus $\gamma = 0.1$. Higher γ improves performance in the upper part of the space where the confirmation bias is strongest. **c)** Optimizing for performance, the optimal γ^* depends on the task. Where the confirmation bias had been strongest, optimal performance is achieved with a stronger leak term. **d)** Model performance when the optimal γ^* from (c) is used in each task. **e)** Comparing the ideal observer to (d), the ideal observer still outperforms the model but only in the upper part of the space. **f)** Temporal weight slopes when using the optimal γ^* are flat everywhere. The models reproduce the change in slopes seen in the data only when γ is fixed across tasks (compare Figure S5).

289 across the space of category and sensory information (Figure 5a-b). We found that in the LSHC
290 regime where the confirmation bias had been strongest, the larger value of γ counteracts the bias
291 and leads to better performance, but in the HSLC regime where there had been no confirmation
292 bias, the optimal γ is zero (Figure 5c). We thus see that the optimal value of γ depends on the
293 task statistics, i.e. the balance of sensory information and category information: the stronger the
294 primacy effect or confirmation bias measured above, the higher γ must be to correct for it (Figure
295 5d). Analogous results were found for the variational model (Supplemental Figure S6).

296 We next asked what the effect would be on the model’s temporal weights if it could utilize the
297 best γ for each task. We found that the γ -optimized model displayed near-flat weights across the
298 entire space of tasks (Figure 5e). Our data therefore imply that either the brain does not optimize
299 its leak to the statistics of the current task, or that it does so on a timescale that is slower than a
300 single experimental session (roughly 1hr, Methods).

301 Predictions for Neurophysiology

302 Both the sampling and variational models induce a confirmation bias by creating an “attractor”
303 dynamic between different levels of the cortical hierarchy – the decision-making area and the relevant
304 sensory areas. Our model therefore makes a number of novel and testable neurophysiological
305 predictions.

306 First, our model predicts that both “choice probabilities” (Britten et al., 1996; Cumming and
307 Nienborg, 2016) and “differential correlations” (Moreno-Bote et al., 2014) in populations of task-
308 relevant sensory neurons will be stronger in contexts where category information is high and sensory
309 information is low, i.e. when subjects exhibit primacy effects (Wimmer et al., 2015; Haefner et al.,
310 2016). This is because the feedback from the decision-making to sensory areas in our model ex-
311 plicitly biases the sensory representation *in the direction that encodes the stimulus strength*, which
312 is the f' -direction (Tajima et al., 2016; Lange and Haefner, 2020). Our model is thus consistent
313 with recent evidence that noise correlations contain a task-dependent component in the f' direction
314 (Bondy et al., 2018).

315 Second, our model predicts that apparent attractor-dynamics measured in both sensory and
316 decision-making areas are in fact driven by inter- rather than within-area dynamics, and will de-
317 pend on the decision-making context. In particular, categorization tasks should induce a stronger
318 confirmation bias, and hence stronger attractor-like dynamics, than equivalent estimation tasks,
319 as was recently reported (Tajima et al., 2017). This observation, as well as our above prediction,
320 contrasts with classic attractor models which posit a recurrent feedback loop *within* a decision
321 making area (Wang, 2008; Wimmer et al., 2015).

322 Discussion

323 Our work makes three main contributions. First, we show that online inference in a hierarchical
324 model can result in characteristic task-dependent temporal biases, and further that such biases
325 naturally arise in two specific families of biologically-plausible approximate inference algorithms.
326 Second, explicitly modeling the mediating sensory representation allows us to partition the infor-
327 mation in the stimulus about the category into two parts – “sensory information” and “category
328 information” – defining a novel two-dimensional space of possible tasks. Third, we collect new data
329 confirming a critical prediction of our theory, namely that individual subjects’ temporal biases
330 change depending on the nature of the information in the stimulus. These results strongly suggest
331 that the discrepancy in temporal biases reported by previous studies is resolved by considering how
332 their tasks trade off sensory and category information.

333 The “confirmation bias” emerges in our models as the result of four key assumptions. Our first
334 assumption is that inference in evidence integration tasks is hierarchical, and that the brain ap-
335 proximates the posterior distribution over both the category, C , and intermediate sensory variables,
336 x . This is in line with converging evidence that populations of sensory neurons encode posterior
337 distributions of corresponding sensory variables (Lee and Mumford, 2003; Yuille and Kersten, 2006;
338 Berkes et al., 2011; Beck et al., 2013) incorporating dynamic prior beliefs via feedback connections
339 (Lee and Mumford, 2003; Yuille and Kersten, 2006; Beck et al., 2013; Nienborg and Roelfsema,
340 2015; Tajima et al., 2016, 2017; Orbán et al., 2016; Haefner et al., 2016; Lange and Haefner, 2020),
341 which contrasts with other probabilistic theories in which only the likelihood is represented in
342 sensory areas (Ma et al., 2006; Beck et al., 2008; Orhan and Ma, 2017; Walker et al., 2019).

343 Our second key assumption is that evidence is accumulated online. In our models, the belief
344 over C is updated based only on the posterior from the previous step and the current posterior over
345 x . This can be thought of as an assumption that the brain does not have a mechanism to store
346 and retrieve earlier frames veridically, but must make use of currently available summary statistics.
347 This is consistent with drift-diffusion models of decision-making (Gold and Shadlen, 2007). As
348 mentioned in the main text, the assumptions until now – hierarchical inference with online updates
349 – do not entail any temporal biases for an ideal observer.

350 Third, we implemented hierarchical online inference making specific assumptions about the
351 limited representational power of sensory areas. In the sampling model, we assumed that the brain
352 can draw a limited number of independent samples of x per update to C . Interestingly, we found
353 that in the small sample regime, the model is inherently unable to account for the prior bias of
354 C on x in its updates to C . Existing neural models of sampling typically assume that samples
355 are distributed temporally (Hoyer and Hyvärinen, 2003; Fiser et al., 2010), but it has also been
356 proposed that the brain could run multiple sampling “chains” distributed spatially (Savin and
357 Denève, 2014). The relevant quantity for our model is the total *effective* number of independent
358 samples that can be generated, stored, and evaluated in a batch to compute each update. The
359 more samples, the smaller the bias predicted by this model.

360 We similarly limited the representational capacity of the variational model by enforcing that the
361 posterior over x is unimodal, and that there is no explicit representation of dependencies between
362 x and C . Importantly, this does not imply that x and C do not influence each other. Rather, the
363 Variational Bayes algorithm expresses these dependencies in the *dynamics* between the two areas:
364 each update that makes $C = +1$ more likely pushes the distribution over x further towards $+1$,
365 and vice versa. Because the number of dependencies between variables grows exponentially, such
366 approximations are necessary in variational inference with many variables (Fiser et al., 2010). The
367 Mean Field Variational Bayes algorithm that we use here has been previously proposed
368 as a candidate algorithm for neural inference (Raju and Pitkow, 2016).

369 The assumptions up to now predict a primacy effect but cannot account for the observed recency
370 effects. When we incorporate a leak term in our models, they reproduce the observed range of biases
371 from primacy to recency. The existence of such a leak term is supported by previous literature
372 (Usher and McClelland, 2001; Bogacz et al., 2006). Further, it is normative in our framework
373 in the sense that reducing the bias in the above models improves performance (Figure 5). The
374 optimal amount of bias correction depends on the task statistics: in the LSHC regime where the
375 confirmation bias is strongest, a higher γ is needed to correct for it. While it is conceivable that
376 the brain would optimize this leak term to the task (Brunton et al., 2013; Piet et al., 2018), our
377 data suggest the leak term is stable across our LSHC and HSLC conditions, or adapted slowly.

378 It has been proposed that post-decision feedback biases subsequent perceptual estimations
379 (Stocker and Simoncelli, 2007; Talluri et al., 2018). While in spirit similar to our confirmation bias
380 model, there are two conceptual differences between these models and our own: First, the feedback

381 from decision area to sensory area in our model is both continuous and online, rather than condi-
382 tioned on a single choice after a decision is made. Second, our models are derived from an ideal
383 observer and only incur bias due to approximations, while previously proposed “self-consistency”
384 biases are not normative and require separate justification.

385 Alternative models have been previously proposed to explain primacy and recency effects in
386 evidence accumulation. Kiani et al. (2008) suggested that an integration-to-bound process is more
387 likely to ignore later evidence even when task-relevant stimuli are of a fixed duration (Kiani et al.,
388 2008). Deneve (2012) showed that simultaneous inference about stimulus strength and choice and
389 in tasks with trials of variable difficulty can lead to either a primacy or a recency effect (Deneve,
390 2012). However, both models of evidence integration are based entirely on total information per
391 frame (i.e. $p(C|e_f)$) and hence cannot explain the difference between the data for the LSHC and
392 the HSLC conditions since both conditions are matched in terms of total information. In general,
393 *any* model based only on $p(C|e_f)$ cannot explain the pattern in our data. While such a model can
394 coexist with the confirmation bias dynamic proposed by our model, it is not sufficient to explain the
395 pattern in our data for which the trade-off between sensory- and category-information is crucial.

396 It has also been proposed that primacy effects could be the result of near-perfect integration
397 of an adapting sensory population (Wimmer et al., 2015; Yates et al., 2017). For this mechanism
398 to explain our full results, however, the sensory population would need to become *less* adapted
399 over frames in our HSLC condition, while at the same time *more* adapted in the LSHC condition.
400 We are unaware of such an adaptation mechanism in the literature. Further, although the circuit
401 dynamics of sensory populations could in principle explain our behavioral results, this would not
402 predict top-down neural effects such as the task-dependence of the dynamics of sensory populations
403 (Tajima et al., 2017) nor the origin and prevalence of differential correlations (Bondy et al., 2018),
404 both of which are consistent with our model, as described above.

405 Models of “leaky” evidence accumulation are known to result in recency effects (Usher and
406 McClelland, 2001; Kiani et al., 2008; Brunton et al., 2013; Glaze et al., 2015). Interestingly, leaky
407 evidence accumulation has also been shown to be optimal in non-stationary environments (Glaze
408 et al., 2015) and could thus in principle indicate that subjects assume such non-stationarity in our
409 HSLC condition. However, this explanation alone cannot explain the presence of primacy effects
410 in the LSHC condition. In sum, while there are numerous existing models that can explain either
411 primacy or recency effects with dedicated mechanisms, ours is the first model to predict the full
412 range of biases and how they may be controlled by the stimulus statistics. Further, because our
413 approximate inference models compute log posterior odds, previously proposed mechanisms like
414 integration to bound are complementary and could be incorporated into our framework.

415 While our focus is on the perceptual domain in which subjects integrate evidence over a timescale
416 on the order of tens or hundreds of milliseconds, analogous principles hold in the cognitive domain
417 over longer timescales. The crucial computational motif underlying our model of the confirmation
418 bias is hierarchical inference over multiple timescales. An agent in such a setting must simultane-
419 ously make accurate judgments of current data (based on the current posterior) and track long-term
420 trends (based on all likelihoods). For instance, Zylberberg et al. (2018) identified an analogous
421 challenge when subjects must simultaneously make categorical decisions each trial (their “fast”
422 timescale) while tracking the stationary statistics of a block of trials (their “slow” timescale), anal-
423 ogous to our LSHC condition. As the authors describe, if subjects base model updates on posteriors
424 rather than likelihoods, they will further entrench existing beliefs (Zylberberg et al., 2018). How-
425 ever, the authors did not investigate order effects; our confirmation bias would predict that subjects’
426 estimates of block statistics is biased towards earlier trials in the block (primacy). Schustek et al.
427 (2018) likewise asked subjects to track information across trials in a cognitive task more analogous
428 to our HSLC condition, and report close to flat weighting of evidence across trials Schustek and

429 Moreno-bote (2018).

430 The strength of the perceptual confirmation bias is directly related to the integration of internal
431 “top-down” beliefs and external “bottom-up” evidence previously implicated in clinical dysfunctions
432 of perception (Jardri and Denève, 2013). Therefore, the differential effect of sensory and category
433 information may be useful in diagnosing clinical conditions that have been hypothesized to be
434 related to abnormal integration of sensory information with internal expectations (Fletcher and
435 Frith, 2009).

436 Hierarchical (approximate) inference on multiple timescales is a common motif across percep-
437 tion, cognition, and machine learning. We suspect that all of these areas will benefit from the
438 insights on the causes of the confirmation bias mechanism that we have described here and how
439 they depend on the statistics of the inputs in a task.

440 Methods

441 Visual Discrimination Task

442 We recruited students at the University of Rochester as subjects in our study. All were compensated
443 for their time, and methods were approved by the Research Subjects Review Board. We found no
444 difference between naive subjects and authors, so all main-text analyses are combined, with data
445 points belonging to authors and naive subjects indicated in Figure 3d.

446 Our stimulus consisted of ten frames of band-pass filtered noise (Beaudot and Mullen, 2006;
447 Nienborg and Cumming, 2014) masked by a soft-edged annulus, leaving a “hole” in the center for
448 a small cross on which subjects fixated. The stimulus subtended 2.6 degrees of visual angle around
449 fixation. Stimuli were presented using Matlab and Psychtoolbox on a 1920x1080px 120 Hz monitor
450 with gamma-corrected luminance (Brainard, 1997). Subjects kept a constant viewing distance of
451 36 inches using a chin-rest. Each trial began with a 200ms “start” cue consisting of a black ring
452 around the location of the upcoming stimulus. Each frame lasted 83.3ms (12 frames per second).
453 The last frame was followed by a single double-contrast noise mask with no orientation energy.
454 Subjects then had a maximum of 1s to respond, or the trial was discarded (Supplemental Figure
455 S1). The stimulus was designed to minimize the effects of small fixational eye movements: (i) small
456 eye movements do not provide more information about either orientation, and (ii) each 83ms frame
457 was too fast for subjects to make multiple fixations on a single frame.

458 The stimulus was constructed from white noise that was then masked by a kernel in the Fourier
459 domain to include energy at a range of orientations and spatial frequencies but random phases
460 (Beaudot and Mullen, 2006; Nienborg and Cumming, 2014; Bondy et al., 2018) (a complete descrip-
461 tion and parameters can be found in the Supplemental Text). We manipulated sensory information
462 by broadening or narrowing the distribution of orientations present in each frame, centered on
463 either $+45^\circ$ or -45° depending on the chosen orientation of each frame. We manipulated category
464 information by changing the proportion of frames that matched the orientation chosen for that
465 trial. The range of spatial frequencies was kept constant for all subjects and in all conditions.

466 Trials were presented in blocks of 100, with typically 8 blocks per session (about 1 hour). Each
467 session consisted of blocks of only HSLC or only LSHC trials (Figure 2). Subjects completed
468 between 1500 and 4400 trials in the LSHC condition, and between 1500 and 3200 trials in the
469 HSLC condition. After each block, subjects were given an optional break and the staircase was
470 reset to $\kappa = 0.8$ and $p_{\text{match}} = 0.9$. p_{match} is defined as the probability that a single frame matched
471 the category for a given trial. In each condition, psychometric curves were fit to the concatenation
472 of all trials from all sessions using the Psignifit Matlab package (Schütt et al., 2016), and temporal
473 weights were fit to all trials below each subject’s threshold.

474 **Low Sensory-, High Category-Information (LSHC) Condition**

475 In the LSHC condition, a continuous 2-to-1 staircase on κ was used to keep subjects near threshold
476 (κ was incremented after each incorrect response, and decremented after two correct responses in
477 a row). p_{match} was fixed to 0.9. On average, subjects had a threshold (defined as 70% correct) of
478 $\kappa = 0.17 \pm 0.07$ (1 standard deviation). Regression of temporal weights was done on all sub-threshold
479 trials, defined per-subject.

480 **High Sensory-, Low Category-Information (HSLC) Condition**

481 In the HSLC condition, the staircase acted on p_{match} while keeping κ fixed at 0.8. Although p_{match}
482 is a continuous parameter, subjects always saw 10 discrete frames, hence the true ratio of frames
483 ranged from 5:5 to 10:0 on any given trial. Subjects were on average $69.5\% \pm 4.7\%$ (1 standard
484 deviation) correct when the ratio of frame types was 6:4, after adjusting for individual biases in the
485 5:5 case. Regression of temporal weights was done on all 6:4 and 5:5 ratio trials for all subjects.

486 **Logistic Regression of Temporal Weights**

487 We constructed a matrix of per-frame signal strengths \mathbf{S} on sub-threshold trials by measuring the
488 empirical signal level in each frame. This was done by taking the dot product of the Fourier-domain
489 energy of each frame as it was displayed on the screen (that is, including the annulus mask applied
490 in pixel space) with a difference of Fourier-domain kernels at $+45^\circ$ and -45° . This gives a scalar
491 value per frame that is positive when the stimulus contained more $+45^\circ$ energy and negative when
492 it contained more -45° energy. Signals were z-scored before performing logistic regression, and
493 weights were normalized to have a mean of 1 after fitting.

494 Temporal weights were first fit using (regularized) logistic regression with different types of
495 regularization. The first regularization method consisted of an AR0 (ridge) prior, and an AR2
496 (curvature penalty) prior. We did not use an AR1 prior to avoid any bias in the slopes, which is
497 central to our analysis.

498 To visualize regularized weights in Figure 3, the ridge and AR2 hyperparameters were chosen
499 using 10-fold cross-validation for each subject, then averaging the optimal hyperparameters across
500 subjects for each task condition. This cross validation procedure was used only for display pur-
501 poses for individual subjects in Figure 3a-c of the main text, while the linear and exponential fits
502 (described below) were used for statistical comparisons. Supplemental Figure S4 shows individual
503 subjects' weights with no regularization.

504 We used two methods to quantify the shape (or slope) of \mathbf{w} : by constraining \mathbf{w} to be either
505 an exponential or linear function of time, but otherwise optimizing the same maximum-likelihood
506 objective as logistic regression. Cross-validation suggests that both of these methods perform sim-
507 ilarly to either unregularized or the regularized logistic regression defined above, with insignificant
508 differences (Supplemental Figure S3). The exponential is defined as

$$\mathbf{w}_f^{\text{exponential}} = \alpha \exp(\beta f) \quad (6)$$

509 where f refers to the frame number. β gives an estimate of the shape of the weights \mathbf{w} over time,
510 while α controls the overall magnitude. $\beta > 0$ corresponds to recency and $\beta < 0$ to primacy. The
511 β parameter is reported for human subjects in Figure 3d, and for the models in Figure 4e,h.

512 The second method to quantify slope was to constrain the weights to be a linear function in
513 time:

$$\mathbf{w}_f^{\text{linear}} = a + \text{slope} \times f \quad (7)$$

514 where $slope > 0$ corresponds to recency and $slope < 0$ to primacy.

515 Figure 3d shows the median exponential shape parameter (β) after bootstrapped resampling of
516 trials 500 times for each subject. Both the exponential and linear weights give comparable results
517 (Supplemental Figure S2).

518 To compute the combined temporal weights across all subjects (in Figure 3a-c), we first esti-
519 mated the mean and variance of the weights for each subject by bootstrap-resampling of the data
520 500 times without regularization. The combined weights were computed as a weighted average
521 across subjects at each frame, weighted by the inverse variance estimated by bootstrapping.

522 Because we are not explicitly interested in the magnitude of \mathbf{w} but rather its *shape* over stimulus
523 frames, we always plot a “normalized” weight, $\mathbf{w}/\text{mean}(\mathbf{w})$, both for our experimental results
524 (Figure 3a-c) and for the model (Figure 4d,g).

525 Approximate inference models

526 We model evidence integration as Bayesian inference in a three-variable generative model (Figure
527 4a) that distills the key features of online evidence integration in a hierarchical model (Haefner
528 et al., 2016). The variables in the model are mapped onto the sensory periphery (e), sensory cortex
529 (x), and a decision-making area (C) in the brain.

530 In the generative direction, on each trial, the binary value of the correct choice $C \in \{-1, +1\}$
531 is drawn from a 50/50 prior. x_f is then drawn from a mixture of two Gaussians:

$$x_f^{(gen)} \sim \begin{cases} \mathcal{N}(+C, \sigma_x^2) & \text{with prob. equal to category info.} \\ \mathcal{N}(-C, \sigma_x^2) & \text{otherwise} \end{cases} \quad (8)$$

532 Finally, each e_f is drawn from a Gaussian around x_f :

$$e_f^{(gen)} \sim \mathcal{N}(x_f, \sigma_e^2) \quad (9)$$

533 When we model inference in this model, we assume that the subject has learned the correct model
534 parameters, even as parameters change between the two different conditions. This is why we ran
535 our subjects in blocks of only LSHC or HSLC trials on a given day.

536 Category information in this model can be quantified by the probability that $x_f^{(gen)}$ is drawn
537 from the mode that matches C . We quantify sensory information as the probability with which an
538 ideal observer can recover the sign of x_f . That is, in our model sensory information is equivalent
539 to the area under the ROC curve for two univariate Gaussian distributions separated by a distance
540 of 2, which is given by

$$\text{sensory info.} = \Phi(\sqrt{2}/\sigma_e) \quad (10)$$

541 where Φ is the inverse cumulative normal distribution.

542 Because the effective time per update in the brain is likely faster than our 83ms stimulus frames,
543 we included an additional parameter n_U for the number of online belief updates per stimulus frame.
544 In the sampling model described below, we amortize the per-frame updates over n_U steps, updating
545 n_U times per frame using $\frac{1}{n_U} \text{LL}O_f$. In the variational model, we interpret n_U as the number of
546 coordinate ascent steps.

547 Simulations of both models were done with 10000 trials per task type and 10 frames per trial.
548 To quantify the evidence-weighting of each model, we used the same logistic regression procedure
549 that was used to analyze human subjects’ behavior. In particular, temporal weights in the model
550 are best described by the exponential weights (equation (6)), so we use β to characterize the model’s
551 biases.

552 Sampling model

553 The sampling model estimates $p(e_f|C)$ using importance sampling of x , where each sample is
 554 drawn from a pseudo-posterior using the current running estimate of $p_{f-1}(C) \equiv p(C|e_1, \dots, e_{f-1})$ as
 555 a marginal prior:

$$x_f^{(s)} \sim Q(x) \propto p(e_f|x_f) \sum_c p(x_f|C=c) p_{f-1}(C=c) \quad (11)$$

556 Using this distribution, we obtain the following unnormalized importance weights.

$$\hat{w}^{(s)} = \left(\sum_c p(x_f^{(s)}|C=c) p_{f-1}(C=c) \right)^{-1} \quad (12)$$

In the self-normalized importance sampling algorithm these weights are then normalized as follows,

$$\hat{w}^{(s)} = \frac{w^{(s)}}{\sum_i w^{(i)}},$$

557 though we found that this had no qualitative effect on the model’s ability to reproduce the trends
 558 in the data. The above equations yield the following estimate for the log-likelihood ratio needed
 559 for the belief update rule in equation (5):

$$\hat{\text{LLO}}_f = \log \frac{\sum_{s=1}^S p(x_f^{(s)}|C=+1) w^{(s)}}{\sum_{s=1}^S p(x_f^{(s)}|C=-1) w^{(s)}} \quad (13)$$

560 In the case of infinitely many samples, these importance weights exactly counteract the bias intro-
 561 duced by sampling from the posterior rather than likelihood, thereby avoiding any double-counting
 562 of the prior, and hence, any confirmation bias. However, in the case of finite samples, S , biased
 563 evidence integration is unavoidable.

564 The full sampling model is given in Supplemental Algorithm S1. Simulations in the main text
 565 were done with $S = 5$, $n_U = 5$, normalized importance weights, and $\gamma = 0$ or $\gamma = 0.1$.

566 Variational model

567 The core assumption of the variational model is that while a decision area approximates the pos-
 568 terior over C and a sensory area approximates the posterior over x , no brain area explicitly rep-
 569 resents posterior dependencies between them. That is, we assume the brain employs a *mean field*
 570 *approximation* to the joint posterior by factorizing $p(C, x_1, \dots, x_F|e_1, \dots, e_F)$ into a product of ap-
 571 proximate marginal distributions $q(C) \prod_{f=1}^F q(x_f)$ and minimizes the Kullback-Leibler divergence
 572 between q and p using a process that can be modeled by the Mean-Field Variational Bayes algorithm
 573 (Murphy, 2012).

574 By restricting the updates to be online (one frame at a time, in order), this model can be seen as
 575 an instance of “Streaming Variational Bayes” (Broderick et al., 2013). That is, the model computes
 576 a sequence of approximate posteriors over C using the same update rule for each frame. We thus
 577 only need to derive the update rules for a single frame and a given prior over C ; this is extended
 578 to multiple frames by re-using the posterior from frame $f - 1$ as the prior on frame f .

579 As in the sampling model, this model is unable to completely discount the added prior over
 580 x . Intuitively, since the mean-field assumption removes explicit correlations between x and C , the

581 model is forced to commit to a marginal posterior in favor of $C = +1$ or $C = -1$ and $x > 0$ or
 582 $x < 0$ after each update, which then biases subsequent judgments of each.

583 To keep conditional distributions in the exponential family (which is only a matter of math-
 584 ematical convenience and has no effect on the ideal observer), we introduce an auxiliary variable
 585 $z_f \in \{-1, +1\}$ that selects which of the two modes x_f is in:

$$z_f = \begin{cases} +1 & \text{with probability equal to category info} \\ -1 & \text{otherwise} \end{cases} \quad (14)$$

586 such that

$$x_f \sim \mathcal{N}(z_f C, \sigma_x^2). \quad (15)$$

587 We then optimize $q(C) \prod_{f=1}^F q(x_f)q(z_f)$.

588 Mean-Field Variational Bayes is a coordinate ascent algorithm on the parameters of each ap-
 589 proximate marginal distribution. To derive the update equations for each step, we begin with the
 590 following (Murphy, 2012):

$$\begin{aligned} \log q(x_f) &\leftarrow \mathbf{E}_{q(C)q(z_f)}[\log p(C, x_f, z_f|e_f)] + \text{const} \\ \log q(z_f) &\leftarrow \mathbf{E}_{q(C)q(x_f)}[\log p(C, x_f, z_f|e_f)] + \text{const} \\ \log q(C) &\leftarrow \mathbf{E}_{q(x_f)q(z_f)}[\log p(C, x_f, z_f|e_f)] + \text{const} \end{aligned} \quad (16)$$

591 After simplifying, the new $q(x_f)$ term is a Gaussian with mean given by equation (17) and constant
 592 variance

$$\mu_{x_f} \leftarrow \frac{\sigma_e^2 \mu_C \mu_{z_f} + \sigma_x^2 e_f}{\sigma_e^2 + \sigma_x^2} \quad (17)$$

593 where μ_C and μ_z are the means of the current estimates of $q(C)$ and $q(z)$.

594 For the update to $q(z_f)$ in terms of log odds of z_f we obtain:

$$\text{LPO}_{z_f} \leftarrow \log \frac{p(z_f = +1)}{p(z_f = -1)} + 2 \frac{\mu_{x_f} \mu_C}{\sigma_e^2 + \sigma_x^2}. \quad (18)$$

595 Similarly, the update to $q(C)$ is given by:

$$\text{LPO}_C \leftarrow \log \frac{p(C = +1)}{p(C = -1)} + 2 \frac{\mu_{x_f} \mu_{z_f}}{\sigma_x^2} \quad (19)$$

596 Note that the first term in equation (19) – the log prior – will be replaced with the log posterior
 597 estimate from the previous frame (see Supplemental Algorithm S2). Comparing equations (19) and
 598 (3), we see that in the variational model, the log likelihood odds estimate is given by

$$\text{L}\hat{\text{L}}\text{O}_f = 2 \frac{\mu_{x_f} \mu_{z_f}}{\sigma_x^2} \quad (20)$$

599 Analogously to the sampling model we assume a number of updates n_U reflecting the speed of
 600 relevant computations in the brain relative to how quickly stimulus frames are presented. Unlike
 601 for the sampling model, naively amortizing the updates implied by equation (20) n_U times results
 602 in a stronger primacy effect than observed in the data, since the Variational Bayes algorithm
 603 naturally has attractor dynamics built in. Allowing for an additional parameter η scaling this
 604 update (corresponding to the step size in Stochastic Variational Inference (Hoffman et al., 2013))
 605 seems biologically plausible because it simply corresponds to a coupling strength in the feed-forward
 606 direction. Decreasing η both reduces the primacy effect and improves the model's performance.
 607 Here we used $\eta = 0.05$ in all simulations based on a qualitative match with the data. The full
 608 variational model is given in Algorithm S2.

609 Acknowledgements

610 This work was supported by NEI/NIH awards R01 EY028811-01 (RMH) and T32 EY007125 (RDL,
611 JLY), as well as an NSF/NRT graduate training grant NSF-1449828 (RDL).

612 Author Contributions

613 Author contributions are shown in the following table, where black = significant contribution, gray
= partial contribution, and white = zero or minimal contribution.

	RL	AC	JB	JY	RH
Experiment Design	black	black	white	black	black
Experiment Code	black	gray	white	white	white
Data Collection	white	black	white	white	white
Data Analysis	black	gray	white	gray	gray
Sampling Model	black	black	white	white	black
Variational Model	black	white	black	white	white
Writing	black	gray	gray	gray	black

614

References

- 615
- 616 William H A Beaudot and Kathy T. Mullen. Orientation discrimination in human vision: Psy-
617 chophysics and modeling. *Vision Research*, 46:26–46, 2006.
- 618 Jeff Beck, Katherine Heller, and Alexandre Pouget. Complex Inference in Neural Circuits with
619 Probabilistic Population Codes and Topic Models. *Advances in Neural Information Processing*
620 *Systems*, 25:3068–3076, 2013.
- 621 Jeffrey M. Beck, Wei Ji Ma, Roozbeh Kiani, Tim Hanks, Anne K. Churchland, Jamie Roitman,
622 Michael N. Shadlen, Peter E. Latham, and Alexandre Pouget. Probabilistic Population Codes
623 for Bayesian Decision Making. *Neuron*, 60(6):1142–1152, 2008.
- 624 Pietro Berkes, Gergo Orbán, Máté Lengyel, and József Fiser. Spontaneous Cortical Activity Reveals
625 Hallmarks of an Optimal Internal Model of the Environment. *Science*, 331(January):83–87, 2011.
- 626 C.M. Bishop. *Pattern Recognition and Machine Learning*. Information science and statistics.
627 Springer (New York), 2006.
- 628 Rafal Bogacz, Eric Brown, Jeff Moehlis, Philip Holmes, and Jonathan D. Cohen. The physics of
629 optimal decision making: A formal analysis of models of performance in two-alternative forced-
630 choice tasks. *Psychological Review*, 113(4):700–765, 2006.
- 631 Adrian G. Bondy, Ralf M. Haefner, and Bruce G. Cumming. Feedback determines the structure of
632 correlated variability in primary visual cortex. *Nature Neuroscience*, 21(4):598–606, 2018.
- 633 D. H. Brainard. The psychophysics toolbox. *Spatial Vision*, 10:433–436, 1997.
- 634 K H Britten, W T Newsome, M N Shadlen, S Celebrini, and J A Movshon. A relationship between
635 behavioral choice and the visual responses of neurons in macaque MT. *Vis Neurosci*, 13(1):
636 87–100, 1996.
- 637 Tamara Broderick, Nicholas Boyd, Andre Wibisono, Ashia C Wilson, and Michael I Jordan. Stream-
638 ing variational bayes. *Advances in Neural Information Processing Systems*, 26:1727–1735, 2013.
- 639 Bingni W Brunton, Matthew M. Botvinick, and Carlos D Brody. Rats and humans can optimally
640 accumulate evidence for decision-making. *Science*, 340(6128):95–8, 2013.
- 641 Chris Cremer, Quaid Morris, and David Duvenaud. Reinterpreting Importance-Weighted Autoen-
642 coders. *arXiv*, pages 1–6, 2017.
- 643 Bruce G. Cumming and Hendrikje Nienborg. Feedforward and feedback sources of choice probability
644 in neural population responses. *Current Opinion in Neurobiology*, 37:126–132, 2016.
- 645 Sophie Deneve. Making Decisions with Unknown Sensory Reliability. *Frontiers in Neuroscience*, 6
646 (June):1–13, 2012.
- 647 Jan Drugowitsch, Valentin Wyart, Anne-Dominique Devauchelle, and Etienne Koechlin. Computa-
648 tional Precision of Mental Inference as Critical Source of Human Choice Suboptimality. *Neuron*,
649 92(6):1398–1411, 2016.
- 650 József József Fiser, Pietro Berkes, Gergo Orbán, and Máté Lengyel. Statistically optimal perception
651 and learning: from behavior to neural representations. *Trends in Cognitive Sciences*, 14(3):
652 119–30, 2010.

- 653 Paul C. Fletcher and Chris D. Frith. Perceiving is believing: A Bayesian approach to explaining
654 the positive symptoms of schizophrenia. *Nature Reviews Neuroscience*, 10:48–58, 2009.
- 655 Samuel J Gershman and Jeffrey M. Beck. Complex Probabilistic Inference: From Cognition to
656 Neural Computation. In Ahmed Moustafa, editor, *Computational Models of Brain and Behavior*,
657 chapter Complex Pr, pages 1–17. Wiley-Blackwell, 2016.
- 658 Christopher M. Glaze, Joseph W. Kable, and Joshua I. Gold. Normative evidence accumulation in
659 unpredictable environments. *eLife*, 4:1–27, 2015.
- 660 Joshua I Gold and Michael N. Shadlen. The neural basis of decision making. *Annual review of*
661 *neuroscience*, 30(30):535–574, 2007.
- 662 Ralf M. Haefner, Pietro Berkes, and Jozsef Fiser. Perceptual Decision-Making as Probabilistic
663 Inference by Neural Sampling. *Neuron*, 90(3):649–660, 2016.
- 664 Matthew D. Hoffman, David M. Blei, Chong Wang, and John Paisley. Stochastic variational
665 inference. *Journal of Machine Learning Research*, 14:1303–1347, 2013.
- 666 P. O. Hoyer and A. Hyvärinen. Interpreting neural response variability as monte carlo sampling of
667 the posterior. *Advances in Neural Information Processing Systems*, 17(1):293–300, 2003.
- 668 Renaud Jardri and Sophie Denève. Circular inferences in schizophrenia. *Brain*, 136(11):3227–3241,
669 2013.
- 670 Roozbeh Kiani, Timothy D Hanks, and Michael N. Shadlen. Bounded integration in parietal cortex
671 underlies decisions even when viewing duration is dictated by the environment. *The Journal of*
672 *Neuroscience*, 28(12):3017–3029, 2008.
- 673 Richard D Lange and Ralf M Haefner. Task-induced neural covariability as a signature of Bayesian
674 learning and inference. *bioRxiv*, 2020.
- 675 Tai Sing Lee and David Mumford. Hierarchical Bayesian inference in the visual cortex. *Journal of*
676 *the Optical Society of America A*, 20(7):1434–1448, 2003.
- 677 Wei Ji Ma, Jeffrey M. Beck, Peter E Latham, and Alexandre Pouget. Bayesian inference with
678 probabilistic population codes. *Nature Neuroscience*, 9(11):1432–1438, 2006.
- 679 Rubén Moreno-Bote, Jeffrey M. Beck, Ingmar Kanitscheider, Xaq Pitkow, Peter Latham, and
680 Alexandre Pouget. Information-limiting correlations. *Nature Neuroscience*, 17(10):1410–1417,
681 2014.
- 682 David Mumford. On the computational architecture of the neocortex. *Biological cybernetics*, 251:
683 241–251, 1992.
- 684 Kevin P. Murphy. *Machine Learning: A Probabilistic Perspective*. The MIT Press, 2012.
- 685 W. T. Newsome and E. B. Pare. A selective impairment of motion perception following lesions of
686 the middle temporal visual area (MT). *The Journal of Neuroscience*, 8(6):2201–2211, 1988.
- 687 Rs Nickerson. Confirmation bias: A ubiquitous phenomenon in many guises. *Review of general*
688 *psychology*, 2(2):175–220, 1998.

- 689 Hendrikje Nienborg and Bruce G. Cumming. Decision-related activity in sensory neurons reflects
690 more than a neuron's causal effect. *Nature*, 459(7243):89–92, 2009.
- 691 Hendrikje Nienborg and Bruce G Cumming. Decision-related activity in sensory neurons may
692 depend on the columnar architecture of cerebral cortex. *The Journal of Neuroscience*, 34(10):
693 3579–85, 2014.
- 694 Hendrikje Nienborg and Pieter R. Roelfsema. Belief states as a framework to explain extra-retinal
695 influences in visual cortex. *Current opinion in neurobiology*, 32:45–52, 2015.
- 696 Bruno a Olshausen and D J Field. Sparse coding with an incomplete basis set: a strategy employed
697 by \protect{V1}, 1997.
- 698 Gergő Orbán, Pietro Berkes, József Fiser, and Máté Lengyel. Neural Variability and Sampling-
699 Based Probabilistic Representations in the Visual Cortex. *Neuron*, 92(2):530–543, 2016.
- 700 A. Emin Orhan and Wei Ji Ma. Efficient probabilistic inference in generic neural networks trained
701 with non-probabilistic feedback. *Nature Communications*, 8(138), 2017.
- 702 Art B. Owen. Importance Sampling. In *Monte Carlo theory, methods and examples*, chapter 9.
703 2013.
- 704 Alex T Piet, Ahmed El Hady, and Carlos D. Brody. Rats adopt the optimal timescale for evidence
705 integration in a dynamic environment. *Nature Communications*, 9, 2018.
- 706 Alexandre Pouget, Jeffrey M. Beck, Wei Ji Ma, and Peter E Latham. Probabilistic brains: knowns
707 and unknowns. *Nature Neuroscience*, 16(9):1170–8, 2013.
- 708 Rajkumar Vasudeva Raju and Xaq Pitkow. Inference by Reparameterization in Neural Population
709 Codes. *Advances in Neural Information Processing Systems*, 30, 2016.
- 710 David Raposo, Matthew T Kaufman, and Anne K Churchland. A category-free neural population
711 supports evolving demands during decision-making. *Nature Neuroscience*, 17(12):1784–1792,
712 2014.
- 713 Cristina Savin and Sophie Denève. Spatio-temporal representations of uncertainty in spiking neural
714 networks. *Advances in Neural Information Processing Systems*, pages 1–9, 2014.
- 715 Philipp Schustek and Rubén Moreno-bote. Human confidence judgments reflect reliability-based
716 hierarchical integration of contextual information. *bioRxiv*, 2018.
- 717 Heiko H. Schütt, Stefan Harmeling, Jakob H. Macke, and Felix A. Wichmann. Painfree and accu-
718 rate Bayesian estimation of psychometric functions for (potentially) overdispersed data. *Vision
719 Research*, 122:105–123, 2016.
- 720 L Shi and Tl Griffiths. Neural implementation of hierarchical Bayesian inference by importance
721 sampling. *Advances in Neural Information Processing Systems*, pages 1–9, 2009.
- 722 Alan A Stocker and Eero P Simoncelli. A Bayesian Model of Conditioned Perception. *Advances in
723 Neural Infromation Processing Systems*, 2007:1409–1416, 2007.
- 724 Chihiro I. Tajima, Satoshi Tajima, Kowa Koida, Hidehiko Komatsu, Kazuyuki Aihara, and
725 Hideyuki Suzuki. Population code dynamics in categorical perception. *Nature Scientific Re-
726 ports*, 5(August 2015):1–13, 2016.

- 727 Satohiro Tajima, Kowa Koida, Chihiro I. Tajima, Hideyuki Suzuki, Kazuyuki Aihara, and Hidehiko
728 Komatsu. Task-dependent recurrent dynamics in visual cortex. *eLife*, 6:1–27, 2017.
- 729 Bharath Chandra Talluri, Anne E Urai, Konstantinos Tsetsos, Marius Usher, Tobias H Donner,
730 Bharath Chandra Talluri, Anne E Urai, Konstantinos Tsetsos, Marius Usher, and Tobias H
731 Donner. Confirmation Bias through Selective Overweighting of Choice-Consistent Evidence Re-
732 port Confirmation Bias through Selective Overweighting of Choice-Consistent Evidence. *Current*
733 *Biology*, pages 1–8, 2018.
- 734 Marius Usher and James L. McClelland. The Time Course of Perceptual Choice: The Leaky,
735 Competing Accumulator Model. *Psychological Review*, 108(2):550–592, 2001.
- 736 A. Wald and J. Wolfowitz. Optimum Character of the Sequential Probability Ratio Test. *The*
737 *Annals of Mathematical Statistics*, 19(3):326–339, 1948.
- 738 Edgar Y Walker, R. James Cotton, Wei Ji Ma, and Andreas S Tolias. A neural basis of probabilistic
739 computation in visual cortex. *Nature Neuroscience*, 23:122–129, 2019.
- 740 Xiao Jing Wang. Decision Making in Recurrent Neuronal Circuits. *Neuron*, 60(2):215–234, 2008.
- 741 Klaus Wimmer, Albert Compte, Alex Roxin, Diogo Peixoto, Alfonso Renart, and Jaime De Rocha.
742 Sensory integration dynamics in a hierarchical network explains choice probabilities in cortical
743 area MT. *Nature Communications*, 6(6177):1–13, 2015.
- 744 Valentin Wyart, Vincent De Gardelle, Jacqueline Scholl, and Christopher Summerfield. Rhythmic
745 Fluctuations in Evidence Accumulation during Decision Making in the Human Brain. *Neuron*,
746 76(4):847–858, 2012.
- 747 Jacob L. Yates, Il Memming Park, Leor N. Katz, Jonathan W. Pillow, and Alexander C. Huk. Func-
748 tional dissection of signal and noise in MT and LIP during decision-making. *Nature neuroscience*,
749 20(9):1285–1292, 2017.
- 750 Alan Yuille and Daniel Kersten. Vision as Bayesian inference: analysis by synthesis? *Trends in*
751 *Cognitive Sciences*, 10(7):301–308, 2006.
- 752 Ariel Zylberberg, Daniel M Wolpert, and Michael N Shadlen. Counterfactual reasoning underlies
753 the learning of priors in decision making. *Neuron*, 99:1–15, 2018.

Supplemental Information: A confirmation bias in perceptual decision-making due to hierarchical approximate inference

Richard D. Lange^{1,2,*}, Ankani Chatteraj¹,
Jeffrey M. Beck³, Jacob L. Yates¹, Ralf M. Haefner^{1,*}

¹Brain and Cognitive Sciences, University of Rochester, Rochester, NY 14627, USA.

²Computer Science, University of Rochester, Rochester, NY 14627, USA.

³Department of Neurobiology, Duke University, Durham, NC 27708, USA.

*Corresponding authors: rlange@ur.rochester.edu, rhaefne2@ur.rochester.edu.

January 20, 2020

Sensory Information and Category Information in Previous Literature

In this section we justify our categorization of previous studies' stimuli into the low-sensory/high-category information (LSHC) or high-sensory/low-category information (HSLC) regime in relation to Figure 1 and Table S1. While category information and sensory information are well defined in our model, in the brain they will depend on the nature of the intermediate variable x relative to e and C , and those relationships depend on the sensory system under consideration. For instance, a high spatial frequency grating may contain high sensory information to a primate, but low sensory information to a species with lower acuity. Similarly, when "frames" are presented quickly, they may be temporally integrated with the effect of both reducing sensory information and increasing category information. Therefore, the placement of each study in the sensory vs category information space is our best estimate, and we generally only distinguish between high and low along each dimension. Note that for the orientation discrimination task that we designed, we report the *within-subject change* in weights from one task condition to the other, which overcomes the difficulties described above: while we cannot estimate the absolute values of sensory and category information due to our limited knowledge about the nature of the human sensory system's representation even in our task, our two-staircase task design acting on the two kinds of information separately guarantees that there will be a change in both sensory information and category information between the LSHC and HSLC conditions while performance is kept constant.

Studies finding a primacy effect

Kiani et al. (2008) studied the classic motion direction discrimination task in which a monkey views a dynamic random dot motion stimulus with a certain percentage of “coherent” dots moving together and the rest moving randomly (Kiani et al., 2008; Newsome and Pare, 1988). Monkeys were trained to categorize the direction of motion as predominantly leftward or rightward. Since the direction of the coherently moving dots (the signal) does not change over time within a trial, this stimulus contains high category information. Since the motion direction is difficult to perceive for any motion frame, it contains low sensory information (Kiani et al., 2008).

Nienborg et al. (2009) developed a task in which subjects viewed a disc with varying binocular disparity. The disc moved back and forth relative to a reference plane (the surrounding ring), changing every 10ms, at a rate too high for the macaques’ (and humans’) binocular system to resolve, resulting in a percept of a jittering cloud of dots which was located slightly in front of or behind the surrounding ring and blurred in depth (Nienborg – private communication). After 200 frames presented over 2 seconds, subjects judged whether the center disc was in front or behind the reference plane. Since the location of the perceived dot cloud is relatively stable, but itself uncertain with respect to the reference, this stimulus contains high category and low sensory information (Nienborg and Cumming, 2009).

Studies finding a recency effect or flat weighting

In two similar studies by Wyart et al. (2012) and by Drugowitsch et al. (2016), human participants viewed a sequence of eight clearly visible oriented gratings presented for at least 250ms each. Participants reported whether, on average, the tilt of the eight elements fell closer to the cardinal or diagonal axes. These tasks contain high sensory information since for a subject there is little uncertainty about the orientation of any one grating. However they contain low category information since the orientation of any one grating provides only little information about the correct choice (Wyart et al., 2012; Drugowitsch et al., 2016).

Brunton et al. (2013) studied both a visual task and an auditory task where subjects were trained to indicate whether they saw/heard more flashes/clicks on the left or right side of the midline. These task stimuli contain high sensory information since each flash/click is high contrast/loud – well above subjects’ detection thresholds. However, they contain low category information since each flash/click contains only little information about the correct choice (Brunton et al., 2013).

Stimulus details

The stimulus was constructed from white noise that was then masked by a kernel in the Fourier domain to include energy at a range of orientations and spatial frequencies but random phases (Beaudot and Mullen, 2006; Nienborg and Cumming, 2014; Bondy et al., 2018). The Fourier-domain kernel consisted of a product of two probability density functions (PDFs): a von Mises PDF over orientation, and a Rician PDF over spatial frequency. This

is best expressed using polar coordinates in the Fourier domain:

$$K_{\rho\theta} = \text{vonMises}(\theta; \mu_\theta, \kappa) \text{Rician}(\rho; \mu_\rho, \sigma_\rho)$$

where θ is the angular coordinate and ρ is the spatial frequency coordinate. After transforming back from the Fourier domain to an image, we applied a soft circular aperture with a hole cut out in the center for the fixation cross. The full pixel-space mask is defined by the equation

$$M = \underbrace{\exp(-4\hat{\rho}^2)}_{\text{Gaussian aperture}} \times \underbrace{(1 + \text{erf}(10 \times (\hat{\rho} - \tau_{\text{ap}}/w_{\text{im}})))}_{\text{Center cutout for fixation cross}}$$

where $\hat{\rho}$ is the normalized Euclidean distance to the center of the image ($\hat{\rho} = 0$ at the center, and $\hat{\rho} = \sqrt{2}$ at the corners), and erf is the Error Function. τ_{ap} controlled the width of the central cutout, and w_{im} is the total width of the stimulus. To summarize, each stimulus frame, I , was generated according to

$$I = M \otimes \mathcal{F}^{-1} [\mathcal{F}[\mathcal{W}] \otimes K_{\rho\theta}]$$

where \mathcal{F} is the 2D discrete Fourier transform, \otimes is element-wise multiplication of each pixel, and \mathcal{W} is white noise. Images were displayed using Psychtoolbox on a 1920x1080px 120 Hz monitor with gamma-corrected luminance (Brainard, 1997). Using an 8-bit luminance range (0 to 255), each frame was normalized to $127 \pm c$ where c is a contrast parameter. All stimulus parameters are summarized in table S2.

Algorithms

Algorithm S1 Importance Sampling (IS) model for evidence integration

$LPO \leftarrow \log \frac{p(C=+1)}{p(C=-1)}$ ▷ initialize log posterior odds to log prior odds
for $f = 1$ to F **do**
 for $n = 1$ to n_U **do**
 $p_C \leftarrow (1 + \exp(-LPO))^{-1}$ ▷ current posterior that $C = +1$
 $\hat{p}(x) \leftarrow p_C \mathcal{N}(+1, \sigma_x^2) + (1 - p_C) \mathcal{N}(-1, \sigma_x^2)$ ▷ Mixture of Gaussians prior on x
 $Q(x) \leftarrow \hat{p}(x) p(e_f | x)$
 for $s = 1 \dots S$ **do**
 $x^{(s)} \sim Q(x)$ ▷ sensory sample from current posterior
 $p_+^{(s)} \leftarrow p(x^{(s)} | C = +1)$ ▷ contribution of each sample to $C = +1$ pool
 $p_-^{(s)} \leftarrow p(x^{(s)} | C = -1)$ ▷ contribution of each sample to $C = -1$ pool
 $w^{(s)} \leftarrow (\sum_c p(x^{(s)} | C = c) p_{f-1}(C = c))^{-1}$ ▷ (unnormalized) weight of each sample
 end for
 $w \leftarrow w / \sum_{s'} w^{(s')}$ ▷ (optionally) normalize weights
 $p_+^{tot} \leftarrow \sum_s p_+^{(s)} w^{(s)}$ ▷ aggregate evidence for $C = +1$
 $p_-^{tot} \leftarrow \sum_s p_-^{(s)} w^{(s)}$ ▷ aggregate evidence for $C = -1$
 $L\hat{L}O_f \leftarrow \log p_+^{tot} - \log p_-^{tot}$
 $LPO \leftarrow LPO(1 - \gamma/n_U) + L\hat{L}O_f/n_U$ ▷ equations (13,5) amortized for n_U updates
 end for
end for

Algorithm S2 Variational Bayes (VB) model for evidence integration

$LPO \leftarrow \log \frac{p(C=+1)}{p(C=-1)}$ ▷ initialize to log prior odds
for $f = 1$ to F **do**
 $\mu_{z_f} \leftarrow 2p(z_f = +1) - 1$ ▷ initialize μ_{z_f} to the prior
 for $n = 1$ to n_U **do**
 $\mu_C \leftarrow 2(1 + \exp(-LPO_C))^{-1} - 1$ ▷ convert log-odds to mean of C
 $\mu_{x_f} \leftarrow \frac{\sigma_e^2 \mu_C \mu_{z_f} + \sigma_x^2 e_f}{\sigma_e^2 + \sigma_x^2}$ ▷ equation (17)
 $LPO_{z_f} \leftarrow \log \frac{p(z_f=+1)}{p(z_f=-1)} + 2 \frac{\mu_{x_f} \mu_C}{\sigma_x^2 + \sigma_e^2}$ ▷ equation (18)
 $\mu_{z_f} \leftarrow 2(1 + \exp(-LPO_{z_f}))^{-1} - 1$ ▷ convert log-odds to mean of z_f
 $L\hat{L}O_f \leftarrow \frac{2\mu_{x_f} \mu_{z_f}}{\sigma_x^2}$ ▷ Equation (20)
 $LPO \leftarrow LPO(1 - \gamma/n_U) + \eta L\hat{L}O_f/n_U$ ▷ Equations (5) and (19) amortized for n_U updates with update strength η
 end for
end for

Example study	Justification for placement in task space (Figure 1, color-coded)	Suggested stimulus manipulation to change weighting (color-coded)
Brunton et al. (2013), Raposo et al. (2014)	Each click is perceptually clear but only weakly predictive of which side has the higher rate.	Make clicks softer or embed them in noise and increase difference in rates between left and right side.
Wyart et al. (2012), Drugowitsch et al. (2016)	Orientation of each frame is clear but only weakly predictive of which “deck” the orientations were drawn from.	Decrease contrast of each frame or increase pixel noise and reduce variance of orientations within each deck.
Kiani et al. (2008)	Net motion is weak (low coherence) and constant over a trial.	Increase motion coherence but vary net motion direction across stimulus frames within a trial.
Nienborg et al. (2009)	Percept is of a jittering cloud of dots whose depth is close to fixation point.	Increase the distance between cloud and fixation point in depth; vary distance across stimulus frames at a rate resolvable by depth perception

Table S1: Justification of placement of example prior studies in Figure 1c and description of stimulus manipulations that will move it to the opposite side of the category–sensory–information space. Each manipulation corresponds to a prediction about how temporal weighting of evidence should change from primacy (red) to flat/recency (blue), or vice versa, as a result.

Parameter	Description	Values (Units)
μ_ρ	mean spatial frequency	6.90 (cycles per degree)
σ_ρ	spread of spatial frequency	3.45 (cycles per degree)
κ	(inverse) spread of orientation energy	$0 \leq \kappa \leq 0.8$
c	image contrast	22
τ_{ap}	width of central annulus cutout	25 (pixels) or 0.43 ($^\circ$)
w_{im}	full image width & height	120 (pixels) or 2.08 ($^\circ$)

Table S2: Stimulus parameters.

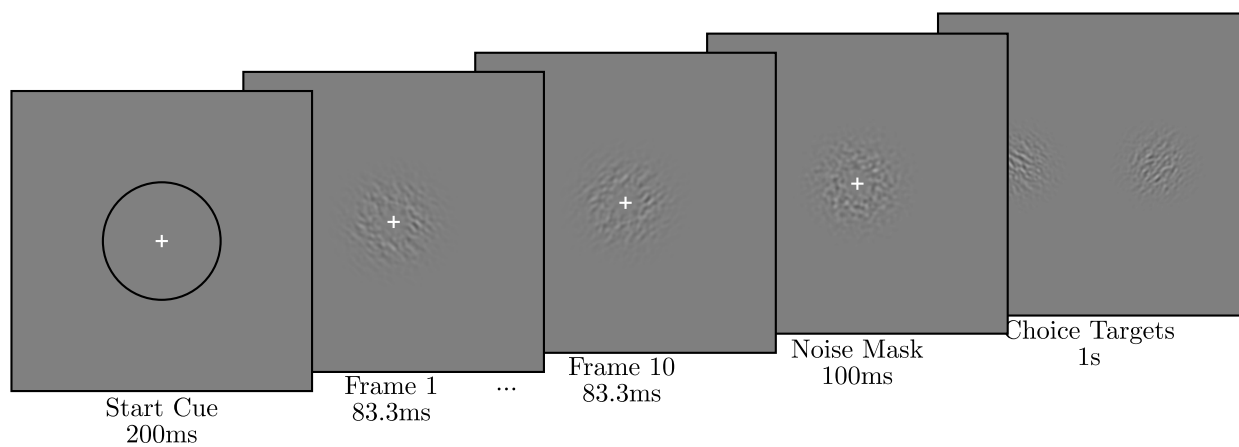


Figure S1: Stimulus timing for each trial in our visual discrimination task

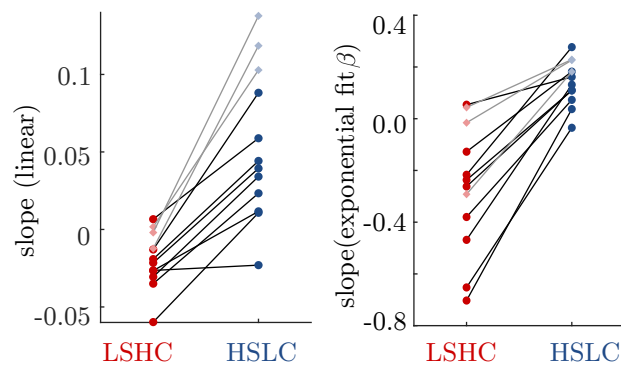


Figure S2: Same as Figure 3d in the main text, comparing slope of \mathbf{w} using a linear fit (left) or an exponential fit (right). Using the linear fit, 11 of 12 subjects individually have a significant increase in slope ($p < 0.05$). Using the exponential fit, 10 of 12 subjects individually have a significant increase in slope ($p < 0.05$).

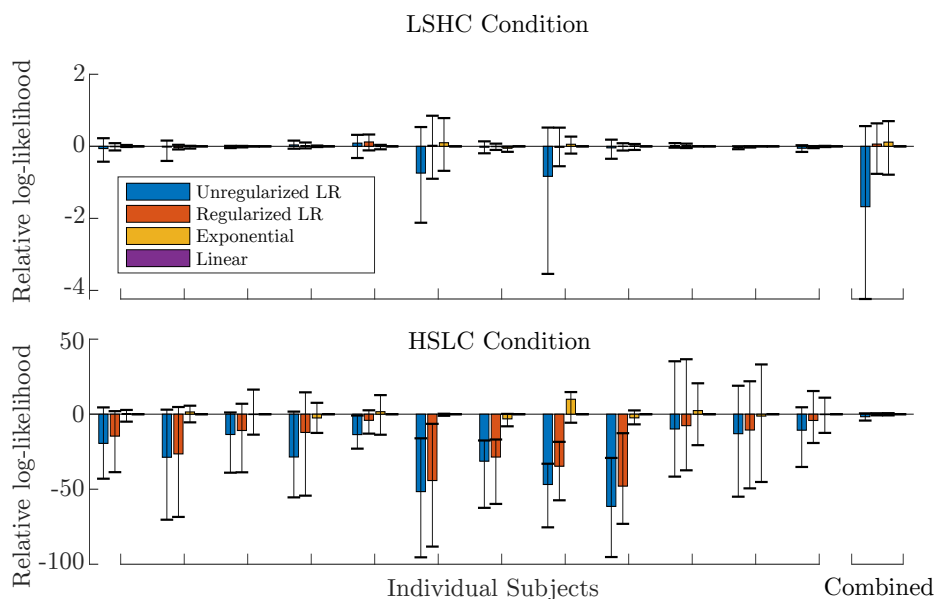


Figure S3: Cross-validation selects linear or exponential shapes for temporal weights, compared to both unregularized and AR2-regularized logistic regression. Panels show 20-fold cross-validation performance of four methods to fit evidence-weighting profiles, separated by task type and by subject. Magnitudes are always relative to the mean log-likelihood of the linear model. Error bars show 50% confidence intervals across folds of shuffled data. “Unregularized LR” refers to standard logistic regression with no regularization. “Regularized LR” refers to the ridge- and AR2-regularized logistic regression objective, where the hyperparameters were chosen to maximize cross-validated fitting performance for each subject. “Exponential” is the 3-parameter model where weights are an exponential function of time (equation (6) plus a bias term). Similarly, the “Linear” model constrains the weights to be a linear function of time as in equation (7), plus a bias term.

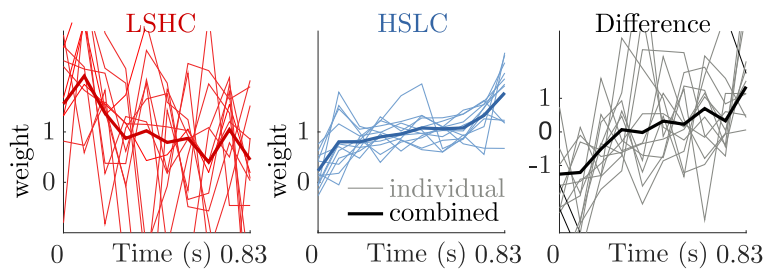


Figure S4: Same as Figure 3a-c in the main text, but with no regularization applied to logistic regression for individual subjects. Both here and in the main text, the “combined” weights are computed using the un-regularized individual weights.

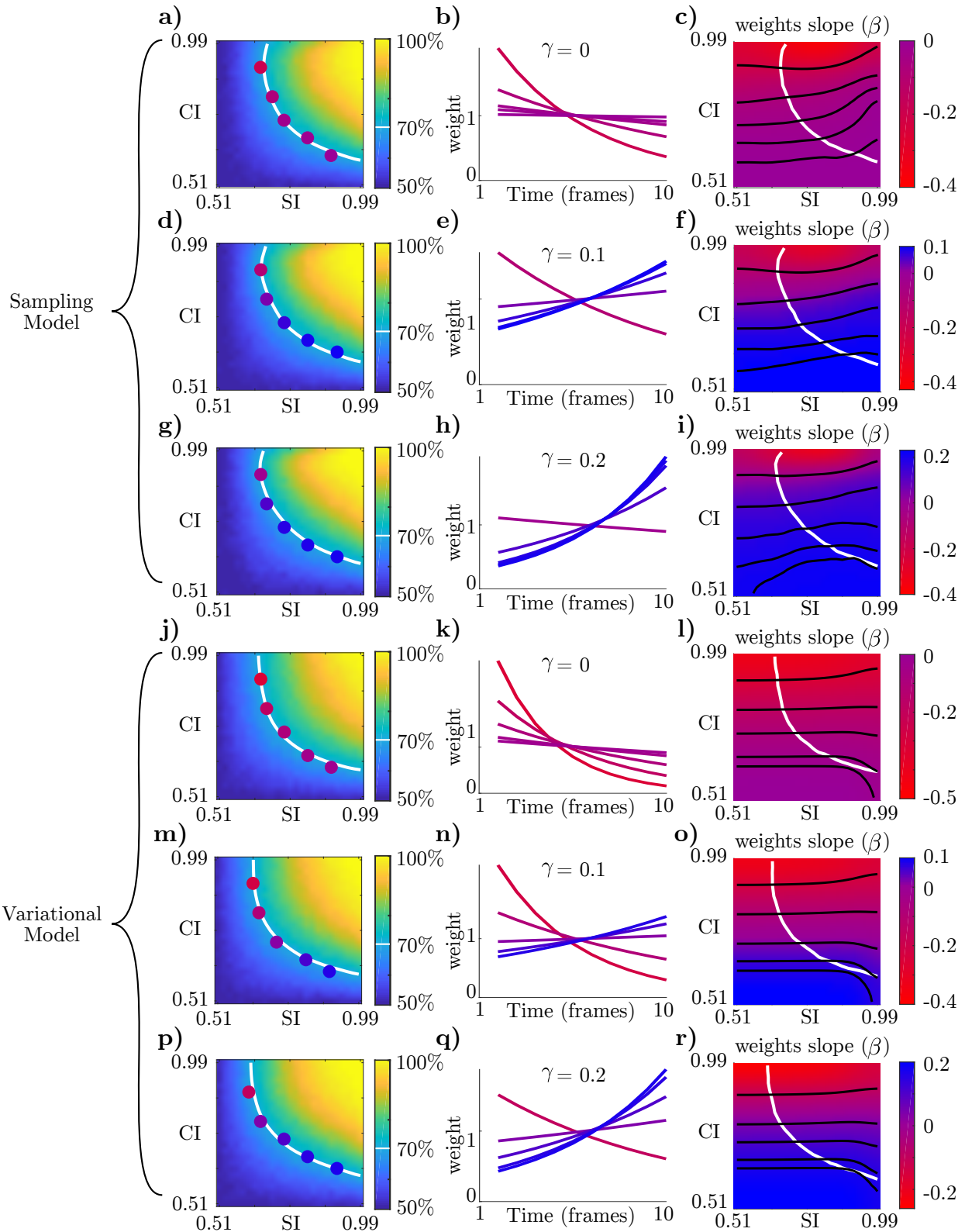


Figure S5: In both models, larger γ increases the prevalence of recency effects across the entire task space. Panels are as in Figure 4 in the main text. **a-c** sampling model with $\gamma = 0$. **d-f** sampling model with $\gamma = 0.1$. **g-i** sampling model with $\gamma = 0.2$. **j-l** variational model with $\gamma = 0$. **m-o** variational model with $\gamma = 0.1$. **p-r** variational model with $\gamma = 0.2$.

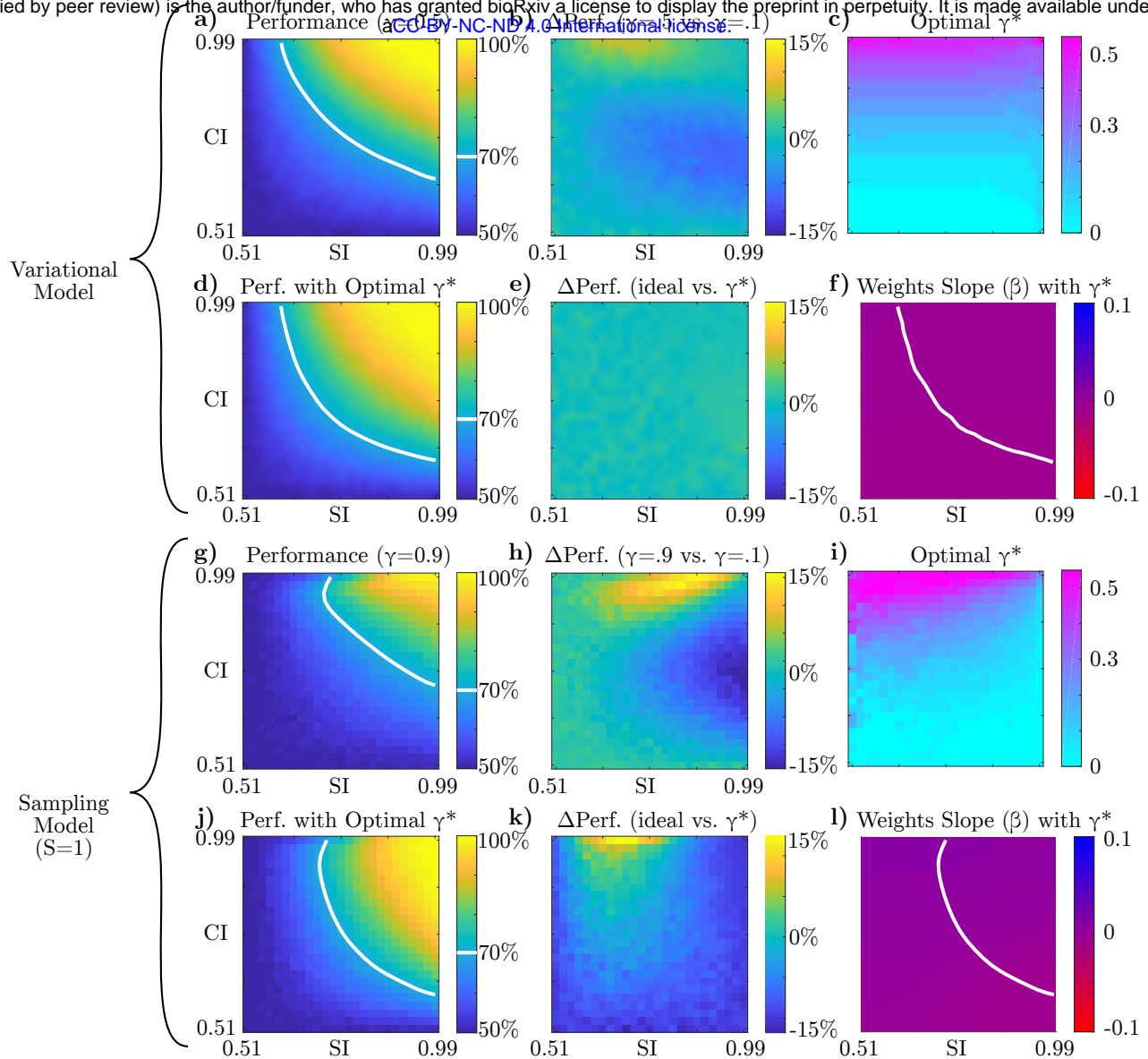


Figure S6: Simulation results for optimal leak (γ) for two further model variations, panels as in Figure 5 in the main text. **a-f** Variational model results. As in the sampling model, we see that the optimal value of γ^* increases with category information, or with the strength of the confirmation bias. **h-l** Sampling model results with $S = 1$ (in the main text we used $S = 5$). Since the sampling model without a leak term approaches the ideal observer in the limit of $S \rightarrow \infty$, the optimal γ^* was close to 0 for much of the space in the main text figure. Here, by comparison, $\gamma^* > 0$ is more common because the $S = 1$ model is more biased.

References

- William H A Beaudot and Kathy T. Mullen. Orientation discrimination in human vision: Psychophysics and modeling. *Vision Research*, 46:26–46, 2006.
- Adrian G. Bondy, Ralf M. Haefner, and Bruce G. Cumming. Feedback determines the

- structure of correlated variability in primary visual cortex. *Nature Neuroscience*, 21(4): 598–606, 2018.
- D. H. Brainard. The psychophysics toolbox. *Spatial Vision*, 10:433–436, 1997.
- Bingni W Brunton, Matthew M. Botvinick, and Carlos D Brody. Rats and humans can optimally accumulate evidence for decision-making. *Science*, 340(6128):95–8, 2013.
- Jan Drugowitsch, Valentin Wyart, Anne-Dominique Devauchelle, and Etienne Koechlin. Computational Precision of Mental Inference as Critical Source of Human Choice Suboptimality. *Neuron*, 92(6):1398–1411, 2016.
- Roозbeh Kiani, Timothy D Hanks, and Michael N. Shadlen. Bounded integration in parietal cortex underlies decisions even when viewing duration is dictated by the environment. *The Journal of Neuroscience*, 28(12):3017–3029, 2008.
- W. T. Newsome and E. B. Pare. A selective impairment of motion perception following lesions of the middle temporal visual area (MT). *The Journal of Neuroscience*, 8(6): 2201–2211, 1988.
- Hendrikje Nienborg and Bruce G. Cumming. Decision-related activity in sensory neurons reflects more than a neuron’s causal effect. *Nature*, 459(7243):89–92, 2009.
- Hendrikje Nienborg and Bruce G Cumming. Decision-related activity in sensory neurons may depend on the columnar architecture of cerebral cortex. *The Journal of Neuroscience*, 34(10):3579–85, 2014.
- Valentin Wyart, Vincent De Gardelle, Jacqueline Scholl, and Christopher Summerfield. Rhythmic Fluctuations in Evidence Accumulation during Decision Making in the Human Brain. *Neuron*, 76(4):847–858, 2012.