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

Richard D. Lange[1,2,*], Ankani Chattoraj[1],
Jeffrey M. Beck[3], Jacob L. Yates[1], Ralf M. Haefner[1,*]

[1]Brain and Cognitive Sciences, University of Rochester, Rochester, NY 14627, USA.
[2]Computer Science, University of Rochester, Rochester, NY 14627, USA.
[3]Department of Neurobiology, Duke University, Durham, NC 27708, USA.
[*]Corresponding authors: rlange@ur.rochester.edu, rhaefne2@ur.rochester.edu.

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

Imagine a doctor trying to infer the cause of a patient's symptoms from an x-ray image. Unsure about the evidence in the image, she asks a radiologist for a second opinion. If she tells the radiologist her suspicion, she may bias his report. If she does not, he may not detect a faint

diagnostic pattern. If the evidence in the image is hard to detect or ambiguous, the radiologist's second opinion, and hence the final diagnosis, may be swayed by the doctors initial hypothesis. We argue that the brain faces a similar problem during perceptual decision-making: any decision-making area combines sequential signals from sensory brain areas, not directly from sensory input. If those signals themselves reflect inferences that combine both prior expectations and sensory evidence, we suggest that this can then lead to an observable confirmation bias (*1*).

Formalizing this idea in the context of approximate Bayesian inference requires extending classic evidence-integration models to include an explicit intermediate sensory representation (Figure 1b). We explicitly model the inferences of the intermediate sensory representation and find that task difficulty is modulated by two distinct types of information: the information between the stimulus and sensory representation (sensory information), and the information between sensory representation and category (category information) (Figure 1b). The balance between these distinct types of information can indeed explain puzzling discrepancies in the literature with regards to the temporal weighting of evidence across a wide range of studies. Even in tasks where all evidence is equally informative about the correct category, existing studies typically report one of three distinct motifs: some find that early evidence is weighted more strongly (a primacy effect) (*2, 3*) some that information is weighted equally over time (as would be optimal) (*4–6*), and some find late evidence being weighted most heavily (a recency effect) (*7*) (Figure 1a,c). There are myriad differences between these studies such as subject species, sensory modality, stimulus parameters, and computational frameworks (*2, 5, 7, 8*). However, none of these aspects can explain their different findings, whereas the differences arise naturally in a hierarchical approximate inference framework.

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
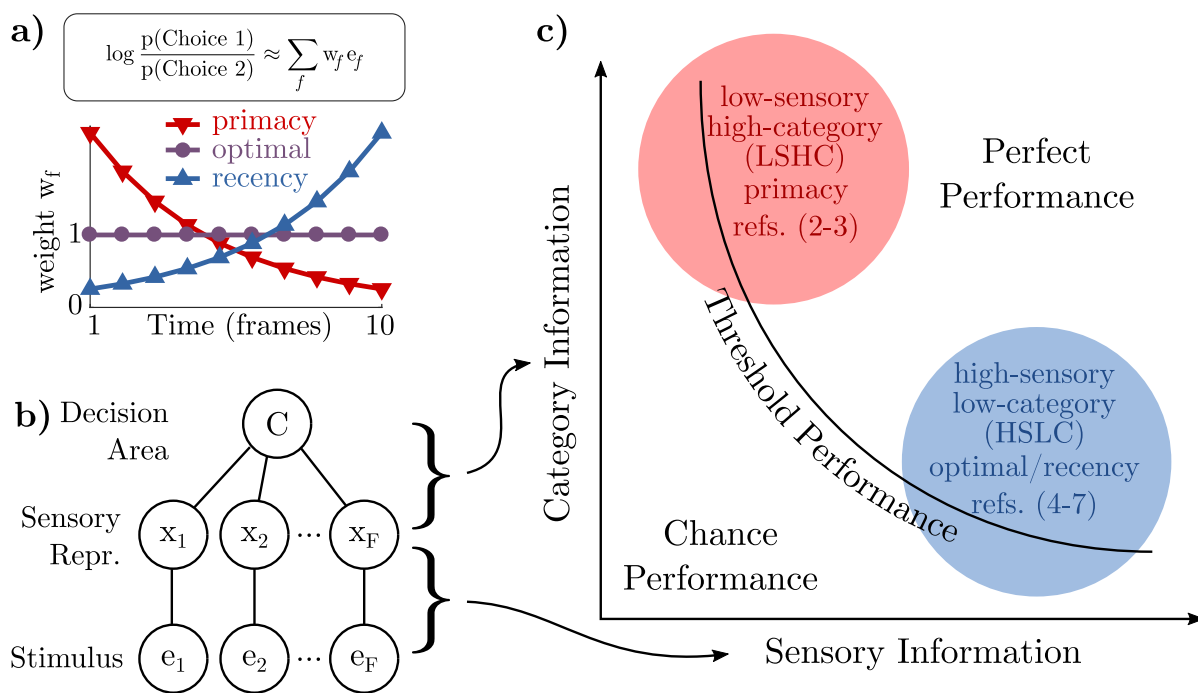
Figure 1: **a)** A subject's "temporal weighting strategy" is an estimate of how their choice is based on a weighted sum of each frame of evidence $e_f$. Three commonly observed motifs are decreasing weights (primacy), constant weights (optimal), or increasing weights (recency). **b)** Uncertainty in the stimulus about the category may be decomposed into uncertainty in each frame about a sensory variable ("sensory information") and uncertainty about the category given the sensory variable ("category information"). **c)** Category information and sensory information may be altered independently, creating a two-dimensional space of possible tasks, where subjects will be at threshold performance whenever the two sources of information are balanced. 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).

3

43 stimulus. On each trial in a typical task, the stimulus consists of multiple "frames" presented in

44 rapid succession. (By "frames" we refer to discrete independent draws of stimulus values that

45 are not necessarily visual). If the evidence in each frame, $e_f$, is independent and the categorical

46 identity of the stimulus is a binary variable $C \in \{-1, +1\}$, then evidence in favor of $C = +1$

47 after $F$ independent frames is $\mathrm{p}(C = +1|e_1, \ldots, e_F) \propto \mathrm{p}(C = +1) \prod_{f=1}^{F} \mathrm{p}(e_f|C = +1)$. The

48 ideal observer reports the most likely category, for instance by reporting the sign of $\log \mathrm{p}(C =$

49 $+1|e_1, \ldots, e_F) - \log \mathrm{p}(C = -1|e_1, \ldots, e_F)$.

50 The ideal observer's performance is limited only by (i) the information about $C$ available on

51 each frame, $p(e_f|C)$, and (ii) the number of frames per trial. In the brain, however, a decision-

52 making area computing a belief about the correct choice only has access to the sensory repre-

53 sentation of the stimulus, which we call $x$, not to the outside stimulus $e$ directly. For example, in

54 a visual task each $e_f$ would be the image on the screen while inferences about $x_f$ are represented

55 by the concurrent activity of relevant neurons in visual cortex. This implies that the information

56 between the stimulus and category can be partitioned into the information between the stimulus

57 and the sensory representation, and the information between sensory representation and cat-

58 egory, which we call "sensory information" and "category information," respectively (Figure

59 1b). These two kinds of information span a two-dimensional space with a task being defined by

60 a single point (Figure 1c).

61 To illustrate this difference, consider the classic dot motion task (*9*) and the Poisson clicks

62 task (*5*), which occupy opposite locations in the space spanned by sensory and category infor-

63 mation. In the classic low-coherence dot motion task, subjects view a cloud of moving dots,

64 some percentage of which move "coherently" in one direction. Here, sensory information is

65 low since evidence about the net motion direction at any time is weak. Category information,

66 on the other hand, is high, since knowing the "true" motion on a single frame would be highly

67 predictive of the correct choice (and of motion on subsequent frames). In the Poisson clicks

4

task, subjects hear a random sequence of clicks in each ear and must report the side with the higher rate. Here, sensory information is high since each click is well above sensory thresholds, but category information is low since knowing the side on which a single click was presented provides only little information about the correct choice (and the side of the other clicks). Another way to think about category information is as "temporal coherence" of the stimulus: the more each frame of evidence is predictive of the correct choice, the more the frames must be predictive of each other, whether a frame consists of visual dots or of auditory clicks. Note that our distinction between sensory and category information is different from the well-studied distinction between internal and external noise; in general, both internal and external noise will reduce the amount of sensory and category information.

If we assume that the sensory representation, which itself is an inference about the actual stimulus, incorporates prior expectations (*10–12*), then, as we show below, approximate inference models predict that this will lead to a primacy effect when sensory information is low and category information is high, but not when sensory information high and category information is low. Indeed, a qualitative placement of prior studies in the space spanned by these two kinds of information demonstrates that studies that find early weighting are located in the upper left quadrant (low-sensory/high-category or LSHC) and studies with equal or late weighting in the lower right quadrant (high-sensory/low-category or HSLC) (Figure 1c). This suggests that the different trade-off between sensory information and category information may indeed underlie differences in temporal weighting seen in previous studies. Further, with this framework it is straightforward to predict how simple changes in stimulus statistics should change the temporal weighting (Table S1). To test this critical model prediction, we designed a visual discrimination task with two stimulus conditions that correspond to the two opposite sides of this task space, while keeping all other aspects of the design the same (Figure 2a). If our theory is correct, then we should be able to change individual subjects' temporal weighting strategy simply by

5

93  changing the sensory–category information trade-off.

94  The stimulus in our task consisted of a sequence of ten visual frames (83ms each). Each

95  frame consisted of band-pass-filtered white noise with excess orientation power either in the

96  $-45°$ or the $+45°$ orientation (*13*) (Figure 2b,d). On each trial, there was a single true orien-

97  tation category, but individual frames might differ. At the end of each trial, subjects reported

98  whether the stimulus was oriented predominantly in the $-45°$ or the $+45°$ orientation. The

99  stimulus was presented as an annulus around the fixation marker in order to minimize the effect

100  of small fixational eye movements (Supplemental Methods).

101  If the brain represents the orientation in each frame, then sensory information in our task is

102  determined by how well each frame determines the orientation of that frame (i.e. the amount of

103  "noise" in each frame), and category information is determined by the probability that any given

104  frame's orientation matches the trial's category. For a ratio of $5 : 5$, a frame's orientation does

105  not predict the correct choice and category information is zero. For a ratio of $10 : 0$, knowledge

106  of the orientation of a single frame is sufficient to determine the correct choice and category

107  information is high. For a more detailed discussion, see Supplementary Text.

108  Using this stimulus, we tested 12 human subjects (9 naive and 3 authors) comparing two

109  conditions intended to probe the difference between the LSHC and HSLC regimes. Starting

110  with both high sensory and high category information, we either ran a staircase lowering the

111  sensory information while keeping category information high, or we ran a staircase lowering

112  category information while keeping sensory information high (Figure 2a). These are the LSHC

113  and HSLC conditions, respectively (Figure 2b,d). For each condition, we used logistic regres-

114  sion to infer, for each subject, the influence of each frame onto their choice. Subjects' overall

115  performance was matched in the two conditions by defining threshold performance as 70%

116  correct (supplementary materials).

117  In agreement with our hypothesis, we find predominantly flat or decreasing temporal weights
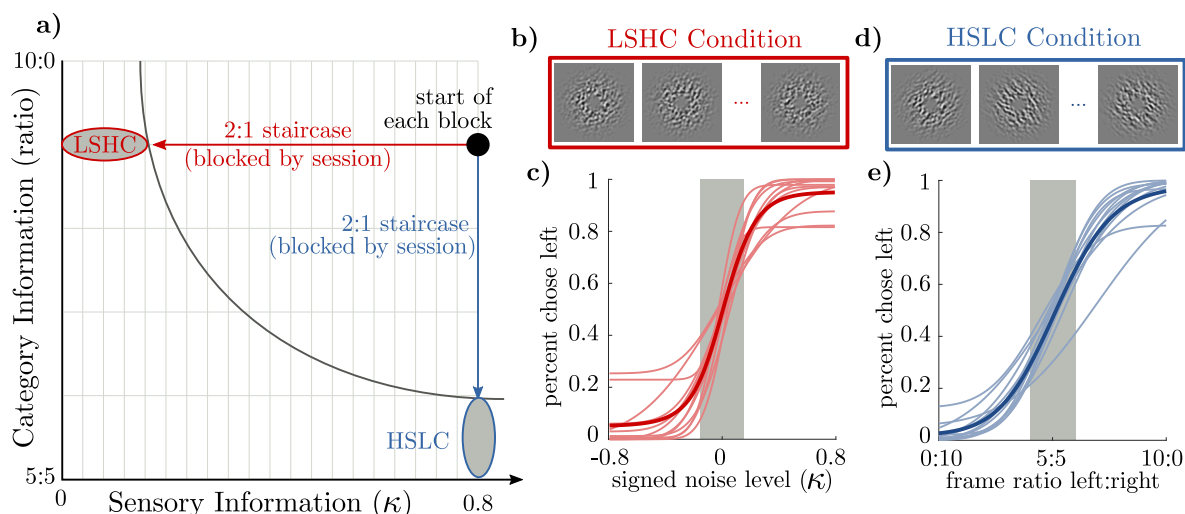
6

Figure 2: Summary of experiment design. **a)** In our task, category information is determined by the ratio of frame categories, and sensory information is determined by a noise parameter $\kappa$. At the start of each block, we reset the staircase to the same point, with category information at $9 : 1$ and $\kappa$ at $0.8$. We then ran a 2-to-1 staircase either on $\kappa$ or on category information, always the same on a given day. The LSHC and HSLC ovals indicate sub-threshold trials; only these trials were used in 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 lines) over the $\kappa$ staircase. Shaded gray area indicates the median threshold level across all subjects. **d)** Visualization of frames in the HSLC condition. Each frame the orientation is clear, but oreintations change from frame to frame. **e)** Psychometric curves over frame ratios, plotted as in (c).
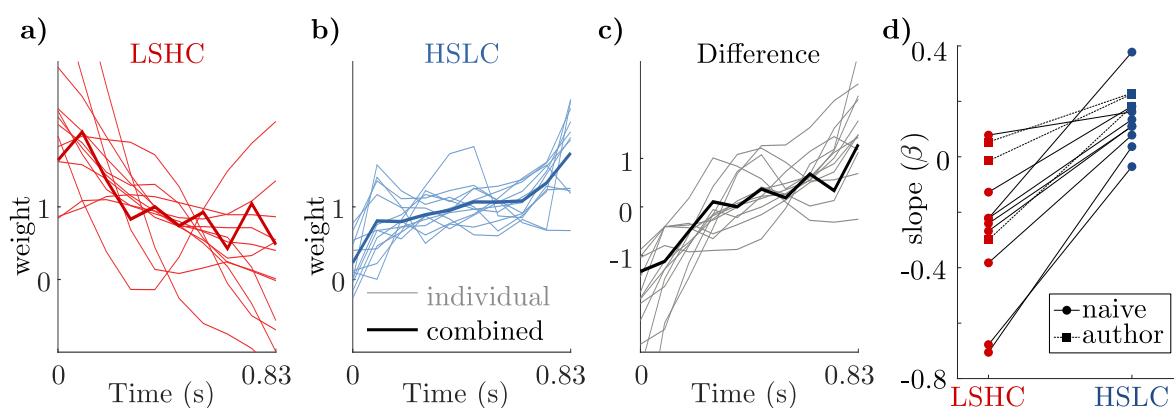
7

Figure 3: Regression of subjects' temporal weights. **a-b)** Temporal weights for individual subjects (thin lines) and the mean across all subjects (thick lines). Weights are always normalized to have a mean of 1. Individual subjects' curves were fit using a cross-validated smoothness term (Supplemental 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. We summarize subjects' temporal weighting strategy with an exponential fit; the slope parameter $\beta > 0$ corresponds to recency and $\beta < 0$ to primacy (Supplemental Methods).

118  when sensory information is low and category information is high (Figure 3a). When the infor-

119  mation is partitioned differently – into high sensory and low categorical information – we find

120  flat or increasing weights (Figure 3b). Despite variability between subjects in each condition, a

121  within-subject comparison revealed that the change in slope between the two conditions was as

122  predicted for every single subject (Figure 2c,d) ($p < 0.05$ for 9 of 12 subjects, bootstrap). This

123  demonstrates that the trade-off between sensory and category information in a task robustly

124  changes subjects' temporal weighting strategy as we predicted, and reconciling the discrepant

125  results in the literature.

126      We will now show that these significant changes in evidence weighting for different stimulus

127  statistics arise naturally in common models of how the brain might implement approximate

128  inference. In particular, we show that both a neural sampling-based approximation (*11, 14–16*)

129  and a parametric (mean-field) approximation (*17,18*) to exact inference can explain the observed

130  pattern of changing temporal weights as a function of stimulus statistics.

131      The crucial assumption in both models is that the brain computes a posterior belief over

132  both $C$ and $x$ given the external evidence, i.e. $\mathrm{p}(x, C|e)$, not just over the categorical variable

133  $C$. This assumption differs from some models of approximate inference in the brain that as-

134  sume populations of sensory neurons strictly encode the *likelihood* of the stimulus (*19*), but is

135  consistent with other models from both sampling and parametric families (*11, 12, 18*).

136      In our models, the brain's belief about $x$ depends both on the external evidence, $e$, via the

137  likelihood, but also on the brain's current belief about $C$, via the prior. For a decision-making

138  area in the brain to update its belief about $C$ based on current sensory responses, it needs to

139  account for, or "subtract out" its influence on those sensory responses. Failure to do so will

140  result in "double-counting" evidence presented early in the trial, inducing a positive feedback

141  loop between the sensory area and the decision making area (Figure 4a). The stronger the

142  decision-making area's belief in a particular choice, the more likely the sensory representation

9

143  of $x$ will concur with that belief through the influence of the prior. We call this feedback loop a

144  "perceptual confirmation bias."

145  Importantly, the strength of this confirmation bias depends on the relative amount of sensory

146  and category information in the stimulus (Figure 4a). It is weakest when the posterior over $x$ is

147  dominated by the likelihood, a case that occurs when the category information is much weaker

148  than the sensory information. Conversely, the feedback loop is strongest when the category

149  information is high compared to the sensory information, as assumed in (*11*) who found a

150  primacy effect in their model.

151  To demonstrate and quantify the intuitions laid out above, we implemented approximate on-

152  line inference (observing a single frame at a time) for a discrimination task using two previously

153  proposed frameworks for how inference might be implemented in neural circuits: neural sam-

154  pling (*11, 14–16*) and Mean Field Variational Inference (*17*) (Figure 4). The central operation

155  in either case is the evaluation of the following ratio (Supplemental Methods)

$$\log \frac{\mathrm{p}(e_f|C=+1)}{\mathrm{p}(e_f|C=-1)} = \log \frac{\int_{x_f} \mathrm{p}(e_f|x_f)\mathrm{p}(x_f|C=+1)\mathrm{d}x_f}{\int_{x_f} \mathrm{p}(e_f|x_f)\mathrm{p}(x_f|C=-1)\mathrm{d}x_f} \tag{1}$$

156  which quantifies how much the brain's belief about $C$ should be changed as the result of the

157  current evidence $e_f$ (*20*). The assumption that information about $C$ is fed back to $x$ is op-

158  erationalized differently in each model, but the effects on the models' behavior are the same

159  (Figure 4, Supplemental Figure S5).

160  The neural sampling hypothesis states that variable neural activity over time can be inter-

161  preted as a sequence of samples from the brain's posterior over $x$. The prior belief about $C$

162  biases the distribution from which samples are generated. The canonical way to compute an

163  expectation with respect to one distribution (the likelihood) using samples from another (the

164  posterior) is 'importance sampling' which weights each sample so as to "subtract out" the prior

165  as described above. While this approach is unbiased in the limit of infinitely many samples,
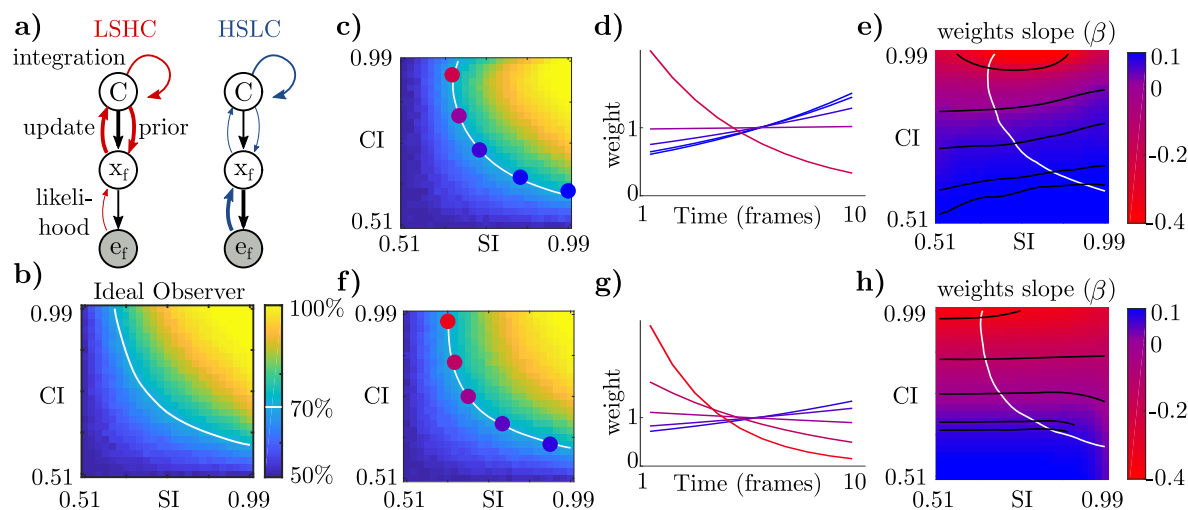
10

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 (inverse width of corresponding conditional distribution). **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 $\gamma = 0.1$.

11

166 it incurs a bias for a finite number – the relevant regime for the brain. The bias is such that it

167 under-corrects for the prior that has been fed back, resulting in a positive feedback loop (see

168 Supplemental Methods). Figure 4b,c shows performance for the ideal observer and for the sam-

169 pling model, respectively, across all combinations of sensory and category information. White

170 lines show threshold performance (70% correct) as in Figure 1c. This model reproduces the

171 primacy effect, and how the temporal weighting changes as the stimulus information changes

172 seen in previous studies. Importantly, it predicted the same within-subject change seen in our

173 data (*11*). However, double-counting the prior alone does not yet explain recency effects (Figure

174 S5a-c,j-l). A simple and biologically-plausible explanation for recency effects is that the brain

175 tries to actively compensate for the prior influence on the sensory representation by subtracting

176 out an estimate of that influence. That is, the brain could do approximate bias correction to

177 mitigate the effect of the confirmation bias. This reduces the primacy effect in the upper left

178 of the task space and leads to a recency effect in the lower right (Figure 4c-h, Supplemental

179 Figure S5), as seen in the data. Interestingly, a linear bias correction term takes the same form

180 as a "leak" term in the classic drift-diffusion framework (Supplemental Methods), and has been

181 shown to be optimal for inference in non-stationary environments (*8, 20*).

182 The second major class of models for how probabilistic inference may be implemented in

183 the brain – based on mean-field parametric representations (*17, 19*) – behaves similarly. These

184 models commonly assume that distributions are encoded *parametrically* in the brain, but that

185 the brain explicitly accounts for dependencies only between subsets of variables, e.g. within the

186 same cortical area. (*18*). We therefore make the assumption that the joint posterior $p(x, C|e)$

187 is approximated in the brain by a product of parametric distributions, $q(x)q(C)$ (*17, 18*). In-

188 ference proceeds by iteratively minimizing the Kullback-Leibler divergence between $q(x)q(C)$

189 and $p(x, C|e)$ (Supplemental Methods). As in the sampling model, the running estimate the

190 category $C$ acts as a prior over $x$. Because this model is unable to explicitly represent posterior

12

191 dependencies between sensory and decision variables, it is biased to commit early either to both

192 $x$ and $C$ being positive or to both $x$ and $C$ being negative. This yields the same behavior as the

193 sampling model: a transition from primacy to flat weights as category information decreases,

194 with recency effects emerging only when approximate bias correction is added (Supplemental

195 Figure 4f-h). Whereas the limited number of samples was the key deviation from optimality

196 in the sampling model, here it is the assumption that the brain represents its beliefs separately

197 about $x$ and $C$ in a factorized form and that its instantaneous belief about $x$ is unimodal.

198 Both models induce a confirmation bias by creating an "attractor" dynamic between differ-

199 ent levels of the cortical hierarchy – the decision-making area and the relevant sensory areas.

200 Our model therefore makes two testable neurophysiological predictions when subjects show

201 a primacy effect: (i) the presence of so-called "differential correlations" (*11, 21*) in popula-

202 tions of task-relevant sensory neurons, and (ii) a reduction of those correlations, as well as any

203 primacy effect, when cortical feedback is inactivated. Our model further predicts that attractor-

204 like dynamics in sensory cortex will depend on the decision-making context, as was recently

205 reported (*22*). This observation, as well as our two novel predictions, contrasts with classic

206 attractor models which posit a recurrent feedback loop within a decision making area (*23*).

207 As in the classic sequential probability ratio test, both models maintain a running estimate

208 of posterior odds over time (*20*). The confirmation bias mechanism is thus complementary to

209 other aspects of evidence integration like "integration to bound" (*2*) or uncertainty over stimulus

210 strength (*24*).

211 In the brain, decisions are not based directly on external evidence but on intermediate repre-

212 sentations . If those intermediate representations themselves in part reflect prior beliefs, and if

213 inference in the brain is approximate, then this is likely to result in a bias. The nature of this bias

214 is directly related to the integration of internal "top-down" beliefs and external "bottom-up" ev-

215 idence previously implicated in clinical dysfunctions of perception (*25*). Importantly, we have

13

shown how the strength of this effect depends on the nature of the information in the task in a way that may generalize to cognitive contexts where the confirmation bias is typically studied.

# References

1. R. Nickerson, *Review of general psychology* **2**, 175 (1998).

2. R. Kiani, T. D. Hanks, M. N. Shadlen, *The Journal of neuroscience* **28**, 3017 (2008).

3. H. Nienborg, B. G. Cumming, *Nature* **459**, 89 (2009).

4. V. Wyart, V. D. Gardelle, J. Scholl, C. Summerfield, *Neuron* **76**, 847 (2012).

5. B. W. Brunton, M. M. Botvinick, C. D. Brody, *Science* **340**, 95 (2013).

6. D. Raposo, M. T. Kaufman, A. K. Churchland, *Nature Neuroscience* **17**, 1784 (2014).

7. J. Drugowitsch, V. Wyart, A.-D. Devauchelle, E. Koechlin, *Neuron* **92**, 1398 (2016).

8. C. M. Glaze, J. W. Kable, J. I. Gold, *eLife* **4**, 1 (2015).

9. W. T. Newsome, E. B. Pare, *The Journal of Neuroscience* **8**, 2201 (1988).

10. T. S. Lee, D. Mumford, *Journal of the Optical Society of America A* **20**, 1434 (2003).

11. R. M. Haefner, P. Berkes, J. Fiser, *Neuron* **90**, 649 (2016).

12. C. I. Tajima, *et al.*, *Nature Scientific Reports* **5**, 1 (2016).

13. W. H. A. Beaudot, K. T. Mullen, *Vision Research* **46**, 26 (2006).

14. P. O. Hoyer, A. Hyvärinen, *Advances in neural information processing systems* **17**, 293 (2003).

15. J. J. Fiser, P. Berkes, G. Orbán, M. Lengyel, *Trends in cognitive sciences* **14**, 119 (2010).

16. G. Orbán, P. Berkes, J. Fiser, M. Lengyel, *Neuron* **92**, 530 (2016).

17. J. Beck, K. Heller, A. Pouget, *Advances in Neural Infromation Processing Systems* **25**, 3068 (2013).

18. R. V. Raju, X. Pitkow, *NIPS* **30** (2016).

19. W. J. Ma, J. M. Beck, P. E. Latham, A. Pouget, *Nature neuroscience* **9**, 1432 (2006).

20. J. I. Gold, M. N. Shadlen, *Annual review of neuroscience* **30**, 535 (2007).

21. R. Moreno-Bote, *et al.*, *Nature Neuroscience* **17**, 1410 (2014).

22. S. Tajima, *et al.*, *eLife* **6**, 1 (2017).

23. X. J. Wang, *Neuron* **60**, 215 (2008).

24. S. Deneve, *Frontiers in Neuroscience* **6**, 1 (2012).

25. R. Jardri, S. Denéve, *Brain* **136**, 3227 (2013).

# Acknowledgments

# Supplementary materials

Materials and Methods

Supplementary Text

Figs. S1-S5

253     Tables S1-S2

254     References *1-15*